



**HAL**  
open science

# On solving parametric polynomial systems and quantifier elimination over the reals: algorithms, complexity and implementations

Huu Phuoc Le

► **To cite this version:**

Huu Phuoc Le. On solving parametric polynomial systems and quantifier elimination over the reals: algorithms, complexity and implementations. Symbolic Computation [cs.SC]. Sorbonne Université, 2021. English. NNT: . tel-03882037v1

**HAL Id: tel-03882037**

**<https://theses.hal.science/tel-03882037v1>**

Submitted on 2 Jan 2022 (v1), last revised 2 Dec 2022 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



**THÈSE DE DOCTORAT DE  
SORBONNE UNIVERSITÉ**

Spécialité

**Informatique**

École Doctorale Informatique, Télécommunications et Électronique (Paris)

Présentée par

**Huu Phuoc LE**

Pour obtenir le grade de

**DOCTEUR de SORBONNE UNIVERSITÉ**

**On solving parametric polynomial systems and quantifier elimination  
over the reals : algorithms, complexity and implementations**

Thèse dirigée par Mohab SAFEY EL DIN

soutenue le vendredi 3 décembre 2021

après avis des **rapporteurs** :

M. Laurent BUSÉ    Research Director, Inria Sophia Antipolis  
M. Éric SCHOST    Professor, University Waterloo

devant le **jury** composé de :

M. Laurent BUSÉ	Research Director, Inria Sophia Antipolis
M. Stef GRAILLAT	Professor, Sorbonne Université
M. Joris VAN DER HOEVEN	Research Director, CNRS
M. Hoon HONG	Professor, North Carolina State University
M. Mohab SAFEY EL DIN	Professor, Sorbonne Université
M. Éric SCHOST	Professor, University Waterloo
M. Bernd STURMFELS	Professor, University of California at Berkeley
Mme. Cynthia VINZANT	Assistant Professor, University of Washington

## Résumé

La résolution de systèmes polynomiaux est un domaine de recherche actif situé entre informatique et mathématiques. Il trouve de nombreuses applications dans divers domaines des sciences de l'ingénieur (robotique, biologie) et du numérique (cryptographie, imagerie, contrôle optimal). Le calcul formel fournit des algorithmes qui permettent de calculer des solutions exactes à ces applications, ce qui pourraient être très délicat pour des algorithmes numériques en raison de la non-linéarité.

La plupart des applications en ingénierie s'intéressent aux solutions réelles. Le développement d'algorithmes permettant de les traiter s'appuie sur les concepts de la géométrie réelle effective; la classe des ensembles semi-algébriques en constituant les objets de base.

Cette thèse se concentre sur trois problèmes ci-dessous, qui apparaissent dans de nombreuses applications et sont largement étudié en calcul formel :

- Classifier les racines réelles d'un système polynomial paramétrique selon les valeurs des paramètres;
- Élimination d'un bloc de quantificateurs;
- Calcul des points isolés d'un ensemble semi-algébrique.

Nous concevons de nouveaux algorithmes symboliques avec une meilleure complexité que l'état de l'art. En pratique, nos implémentations efficaces de ces algorithmes sont capables de résoudre des problèmes hors d'atteinte des logiciels de l'état de l'art.

**Mots-clés.** calcul formel; bases de Gröbner; géométrie algébrique réelle; systèmes polynomiaux paramétriques; élimination de quantificateurs; points isolés réels

## Abstract

Solving polynomial systems is an active research area located between computer sciences and mathematics. It finds many applications in various fields of engineering and sciences (robotics, biology, cryptography, imaging, optimal control). In symbolic computation, one studies and designs efficient algorithms that compute exact solutions to those applications, which could be very delicate for numerical methods because of the non-linearity of the given systems.

Most applications in engineering are interested in the real solutions to the system. The development of algorithms to deal with polynomial systems over the reals is based on the concepts of effective real algebraic geometry in which the class of semi-algebraic sets constitute the main objects.

This thesis focuses on three problems below, which appear in many applications and are widely studied in computer algebra and effective real algebraic geometry:

- Classify the real solutions of a parametric polynomial system according to the values of the parameters;
- One-block quantifier elimination, which is also the computation of the projection of a semi-algebraic set
- Computation of the isolated points of a semi-algebraic set.

We designed new symbolic algorithms with better complexity than the state-of-the-art. In practice, our efficient implementations of these algorithms are capable of solving applications beyond the reach of the state-of-the-art software.

**Keywords.** symbolic computation ; Gröbner bases ; real algebraic geometry ; parametric polynomial systems ; quantifier elimination ; real isolated points

## Acknowledgment

The preparation of this PhD has only been possible thanks to the help and support of many people.

First I would like to thank my advisor Mohab Safey El Din for his supervision and his advices for my research during the past three years. I have learned an incredible amount of things from him. I thank also Jean-Charles Faugère as a supervisor of my master internship, during which I learned the first notions on solving polynomial systems.

I would like to thank Laurent Busé and Éric Schost for their time to read and review my thesis. Their valuable comments improve greatly the quality of the manuscript. In addition, I would like to thank Stef Graillat, Joris van der Hoeven, Hoon Hong, Bernd Sturmfels and Cynthia Vinzant for accepting to be part of the jury of my defense.

I would like to thank Cordian Riener and Stef Graillat for being my mid-term evaluation committee.

I am grateful to have the opportunity to work with Dimitri Manevich, Daniel Plaumann and Timo de Wolff on many interesting problems. From these kind collaborators, I have learned a wide range of scientific knowledge, from theoretical side to applications.

I thank other members of PolSys for their companionship, the scientific discussion and all the helps they willingly provide me all these years. Those are Andrew, Elias, Georgy, Hieu, Irphane, Jérémy, Jocelyn, Jorge, Ludovic, Matias, Olive, Vincent, Rachel, Rémi, Solane and Xuan. In particular, I would like to thank Jérémy for his advices on many things and his help to accomplish my teaching duties.

I would also like to thank Alin Bostan and Frédéric Chyzak with whom I co-organized a joint seminar for PolSys and SpecFun during a year.

Many thanks to all of my teachers without whom I could not be at this stage. Especially to my high school teacher Nguyen Duy Thai Son, he revealed beautiful mathematics to me and inspire me with his passion for solving problems.

My unconditional love is dedicated to my parents who guided me to learn to be a good and happy person. I thank my sister for the birthday presents she has been giving me for the past ten years, unfortunately never on my birthday but only when I came home.

Finally, I thank Minh Ha for accompanying me with a warm heart through this three-year journey. Thank you to her for sharing her time with me, the one who always seems short of time, and for filling our time with joyful memories and a lot of great movies. She is always my infinite source of happiness and courage.

# Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	Motivations and problem statements	5
1.2	State-of-the-art	14
1.2.1	Cylindrical algebraic decomposition	15
1.2.2	One block quantifier elimination	15
1.2.3	Real root classification	18
1.2.4	Computing isolated points in a real algebraic set	19
1.3	Contributions	20
1.3.1	Real root classification	21
1.3.2	One block quantifier elimination	28
1.3.3	Computing the isolated points of a real algebraic set	31
1.4	Organization of the thesis	34
<b>I</b>	<b>Preliminaries</b>	<b>35</b>
<b>2</b>	<b>Basic notions of algebra and geometry</b>	<b>36</b>
2.1	Ideals	36
2.2	Affine algebraic sets	39
2.3	Genericity and changes of variables	42
2.4	Tangent spaces and singularities	44
2.5	Morphisms between affine algebraic sets	45
2.6	Projective algebraic sets	48
2.7	Hilbert series and regular sequences	50
2.8	Cohen-Macaulay rings and determinantal ideals	52
2.9	Noether position and properness	53
<b>3</b>	<b>Gröbner bases</b>	<b>56</b>
3.1	Preliminaries on Gröbner bases	56
3.2	Algebraic elimination using Gröbner bases	60
3.3	Gröbner bases and zero-dimensional ideals	61
3.4	On the computation of Gröbner bases	64
3.4.1	Buchberger algorithm's drawbacks	65
3.4.2	F4/F5 algorithms	65
3.4.3	Complexity issues	65
3.4.4	Change of monomial ordering algorithms	67

<b>4</b>	<b>Basic notions of real algebraic geometry</b>	<b>69</b>
4.1	Real fields . . . . .	69
4.2	Semi-algebraic sets . . . . .	70
4.3	Puiseux series . . . . .	74
4.4	Real root counting . . . . .	76
4.4.1	Notations . . . . .	76
4.4.2	Sturm sequences and Sturm-Habicht sequences . . . . .	76
4.4.3	Hermite quadratic forms . . . . .	82
<b>II</b>	<b>Contributions</b>	<b>86</b>
<b>5</b>	<b>Real root classification algorithms</b>	<b>87</b>
5.1	Introduction . . . . .	88
5.1.1	Problem statement . . . . .	88
5.1.2	Main results . . . . .	90
5.2	Computing sample points in semi-algebraic sets defined by the non-vanishing of polynomials . . . . .	94
5.3	Algorithm based on Sturm-Habicht sequences . . . . .	99
5.4	Parametric Hermite matrices . . . . .	105
5.4.1	Definition . . . . .	105
5.4.2	Gröbner bases and parametric Hermite matrices . . . . .	107
5.4.3	Specialization property of parametric Hermite matrices . . . . .	109
5.4.4	Computing parametric Hermite matrices . . . . .	112
5.5	Algorithms for real root classification . . . . .	115
5.5.1	Algorithm for weak real root classification . . . . .	115
5.5.2	Computing semi-algebraic formulas . . . . .	118
5.6	Complexity analysis . . . . .	121
5.6.1	Degree bounds of parametric Hermite matrices on generic input . . . . .	121
5.6.2	Complexity analysis of our algorithms . . . . .	129
5.7	Practical implementation & Experimental results . . . . .	132
5.7.1	Remark on the implementation of Algorithm 5.4 . . . . .	132
5.7.2	Experiments . . . . .	133
<b>6</b>	<b>One block quantifier elimination for regular polynomial systems of equations</b>	<b>141</b>
6.1	Introduction . . . . .	142
6.1.1	Problem statement . . . . .	142
6.1.2	Main results . . . . .	143
6.2	Algorithm for real root finding . . . . .	145
6.2.1	Safey El Din-Schost algorithm . . . . .	145
6.2.2	Parametric variant of Safey El Din - Schost algorithm . . . . .	146

6.3	One-block quantifier elimination algorithm . . . . .	148
6.3.1	Description . . . . .	148
6.3.2	Correctness . . . . .	151
6.4	Complexity analysis . . . . .	152
6.5	Experiments . . . . .	157
<b>7</b>	<b>Computing totally real hyperplane sections on algebraic curves</b>	<b>160</b>
7.1	Introduction . . . . .	160
7.2	Preliminaries . . . . .	162
7.3	Algorithm for computing totally real hyperplane sections . . . . .	164
7.4	Real algebraic curves in $\mathbb{P}^3$ . . . . .	167
7.5	Plane quartics . . . . .	174
<b>8</b>	<b>Computing the set of isolated points of a real algebraic set</b>	<b>178</b>
8.1	Introduction . . . . .	179
8.1.1	Problem statement . . . . .	179
8.1.2	Main results . . . . .	179
8.2	Candidates for isolated points . . . . .	182
8.2.1	Identification of the candidates . . . . .	182
8.2.2	Computation of candidates . . . . .	184
8.3	The algorithm using roadmaps . . . . .	186
8.3.1	Simplification . . . . .	186
8.3.2	Geometric results . . . . .	187
8.3.3	Description of the algorithm . . . . .	189
8.3.4	Complexity analysis . . . . .	194
8.4	The algorithm of complexity $D^{O(n)}$ . . . . .	195
8.4.1	Geometric results . . . . .	195
8.4.2	Outline of the algorithm . . . . .	199
8.4.3	Computing a value for $e_0$ . . . . .	200
8.4.4	The first variant of <code>Isolated</code> . . . . .	204
8.4.5	Approximations of the candidates . . . . .	206
8.4.6	Complexity analysis . . . . .	209
8.5	Optimizations . . . . .	212
8.5.1	Simple identification of real isolated points . . . . .	212
8.5.2	Limits of critical curves . . . . .	214
8.6	Experimental results . . . . .	216
<b>9</b>	<b>Topics for future research</b>	<b>219</b>
9.1	Resolution of parametric polynomial systems . . . . .	219
9.2	Total real intersection by hyperplanes . . . . .	222
9.3	Computing real isolated points . . . . .	224

# Chapter 1

## Introduction

### 1.1 Motivations and problem statements

This thesis focuses on *polynomial system solving* which is an active research area in between computational mathematics and computer sciences.

As polynomial equations and inequalities allow to model non-linear phenomena, solving polynomial systems finds many applications in several domains, for example, robotics [45, 163, 201], control theory [103, 73], optimization [132, 71], cryptography [121, 62], signal processing [90], biology [186], etc. This problem is intrinsically difficult. For instance, deciding whether a polynomial system has a solution is known to be NP-complete even over a finite field [117].

Polynomial systems are strongly related to algebraic geometry in which one studies *algebraic sets*, i.e., the sets of solutions of polynomial equations over an algebraically closed field such as  $\mathbb{C}$ . Especially when we focus on solutions over  $\mathbb{R}$  (or its generalizations, the real closed fields), the theory of *real algebraic geometry* becomes useful.

In real algebraic geometry, the central objects are *semi-algebraic sets*. A *basic semi-algebraic set* is the set of real solutions of a polynomial system of type

$$f_1 = \dots = f_s = 0, \quad g_1 > 0, \dots, g_r > 0,$$

where  $f_i, g_j$  are polynomials with real coefficients. A semi-algebraic set is a finite union of basic semi-algebraic sets.

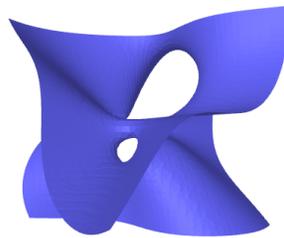


Figure 1.1: The semi-algebraic set defined by the polynomial system

$$x^3y + xz^3 + y^3z + z^3 + 7z^2 + 5z = 0 \wedge x^2 < 25 \wedge y^2 < 25 \wedge z^2 < 25.$$

Various applications in sciences boil down to study properties of semi-algebraic sets. For this purpose, a branch of real algebraic geometry has been developed to design algorithms to investigate semi-algebraic sets with the help of computers; this is also a major topic of our research interest. A few fundamental problems in this direction are

- Computing exactly (at least) one point per connected components of a semi-algebraic set;
- Computing the dimension of semi-algebraic sets;
- Deciding the connectedness between two given points on a semi-algebraic set;
- Computing a description of the projection of a semi-algebraic set.

The design of algorithms for solving these algorithmic problems and many others frequently follows the so-called *critical point method*, a classical technique in optimization and Morse theory. The main principle of this approach is to compute critical points of some well-chosen mapping on the algebraic set under study. The definition of critical points is recalled below.

Let  $V \subset \mathbb{C}^n$  be an algebraic set defined by a sequence  $(f_1, \dots, f_s) \subset \mathbb{C}[x_1, \dots, x_n]$ . Given  $(\varphi_1, \dots, \varphi_m) \subset \mathbb{C}[x_1, \dots, x_n]$ , these polynomials define the mapping

$$\begin{aligned} \varphi : \quad \mathbb{C}^n &\quad \rightarrow \quad \mathbb{C}^m, \\ \mathbf{x} = (x_1, \dots, x_n) &\mapsto (\varphi_1(\mathbf{x}), \dots, \varphi_m(\mathbf{x})). \end{aligned}$$

Under some assumptions on  $V$  (smoothness, equidimensionality,...), the set of critical points of the restriction of  $\varphi$  to  $V$ , denoted by  $\text{crit}(\varphi, V)$  is the simultaneously vanishing set of certain suitable minors of the Jacobian matrix

$$\begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_s}{\partial x_1} & \cdots & \frac{\partial f_s}{\partial x_n} \\ \frac{\partial \varphi_1}{\partial x_1} & \cdots & \frac{\partial \varphi_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial \varphi_m}{\partial x_1} & \cdots & \frac{\partial \varphi_m}{\partial x_n} \end{bmatrix}.$$

In many problems, computations in real algebraic geometry boil down to investigate polynomial systems defining critical points of some well-suited morphisms. Thus, this field of research motivates also the development of efficient algorithmic tools for studying polynomial systems.

Algorithms for polynomial systems are designed following two main paradigms: numerical methods (Newton's method, numerical homotopy continuation method,...) [141, 140, 195, 15] and symbolic methods (multivariate resultants, triangular sets, rational univariate representation, geometric resolution, Gröbner bases,...) [35, 4, 166, 87, 58].

While numerical methods can provide efficiently approximations for the answer, certifying their outputs or guaranteeing the convergence are delicate because of the non-linearity of the

studied systems. By contrast, the outputs of symbolic computations are exact. Many applications in motion planning [34, 128], medical imagery [20, 21], program verification [120, 206, 189] and theorem proving [204, 51] privilege this criterion of exactness. Therefore, in this thesis, we focus on symbolic algorithms for solving polynomial systems. Depending on the nature of input systems or demands of applications, the word “solving” can bear various meanings.

When the considered system has a finite number of solutions in an algebraic closure of the coefficient field, it is called a *zero-dimensional system*. In symbolic computation, enumerating all the solutions of a system usually means computing a representation of the solution set. We will use the following *zero-dimensional parametrization* which historically goes back to Kronecker [124] and is widely used in computer algebra (see, e.g., [40, 78, 33, 102, 131, 82, 86, 139, 87]).

Given a zero-dimensional system in  $\mathbb{Q}[x_1, \dots, x_n]$  whose set of complex solutions is denoted by  $V$ , a zero-dimensional parametrization representing  $V$  consists of

- A *square-free* polynomial  $w \in \mathbb{Q}[u]$  where  $u$  is a new variable;
- A sequence of polynomials  $(v_1, \dots, v_n)$  in  $\mathbb{Q}[u]$  with  $\deg(v_i) < \deg(w)$  such that

$$V = \left\{ \left( \frac{v_1(u)}{w'(u)}, \dots, \frac{v_n(u)}{w'(u)} \right) \mid w(u) = 0 \right\};$$

- A sequence  $(\lambda_1, \dots, \lambda_n) \in \mathbb{Q}^n$  such that

$$u \cdot w' = \sum_{i=1}^n \lambda_i \cdot v_i \pmod{w}.$$

Intuitively,  $u$  coincides with the linear form  $\lambda_1 x_1 + \dots + \lambda_n x_n$  over  $V$ .

**Example 1.1.1.** We consider the zero-dimensional system defined by  $f_1 = f_2 = 0$  (Fig. 1.2) where

$$f_1 = x_1^2 - 2x_2^2 + 2x_1 + 2 \quad \text{and} \quad f_2 = 2x_1x_2 + x_2^2 + x_1 + x_2.$$

A zero-dimensional parametrization for this system consists of  $(\lambda_1, \lambda_2) = (1, 1)$  and

$$\begin{aligned} w &= 7u^4 + 26u^3 + 31u^2 + 8u - 4, \\ v_1 &= -2(8u^3 + 23u^2 + 17u + 2), \\ v_2 &= -2(5u^3 + 8u^2 - 5u - 10). \end{aligned}$$

Many algorithms have been developed to compute zero-dimensional parametrizations among which Gröbner bases [29] and geometric resolutions [87] are two notable approaches.

A Gröbner basis is a finite generating set of an ideal with extra properties which make it a powerful algorithmic tool in computer algebra. They were introduced in Buchberger’s PhD thesis [29] where he also provided the first algorithm to compute them. Even though the worst-case

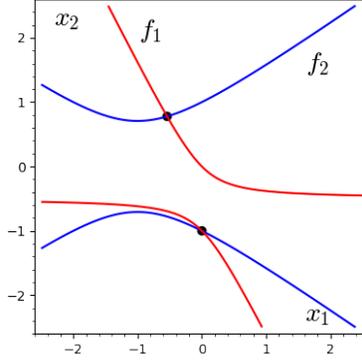


Figure 1.2:  $f_1 = f_2 = 0$ .

complexity of computing Gröbner bases is doubly exponential in the number of variables (see, e.g., [146]), this behavior is reached on only extremely rare systems which are constructed on purpose. In recent decades, new complexity results [133, 80, 81, 130, 131, 7] show that the actual complexity of computing Gröbner is simply exponential in the number of variables for a wide range of systems. Many algorithmic improvements, namely Faugère’s F4/F5 algorithms [59, 60], the FGLM algorithm for ordering change [67] and efficient implementations in the libraries FGB [66], MSOLVE [17] or computer algebra systems (Maple, Magma) reflect the capability of Gröbner bases in solving non-trivial applications. The resolutions of zero-dimensional systems in [166, 58] make use of these advances of Gröbner bases.

The geometric resolution algorithm is a probabilistic algorithm for solving zero-dimensional systems developed by Giusti, Lecerf and Salvy [87]. While the zero-dimensional parametrization has been used since Kronecker [124] by many authors, its computation depends mostly on using Gröbner bases and their related tools. In 1995, Giusti, Heintz, Morais and Pardo [86, 158] rediscovered and improved Kronecker’s approach which does not involve Gröbner bases. Combining these results with the so-called *straight-line program* (see, e.g., [84, 83, 85]), [87] presents the geometric resolution algorithm for solving zero-dimensional systems under some assumptions and estimates its complexity. This algorithm performs an incremental lifting and intersecting process to compute a zero-dimensional parametrization. Its complexity is polynomial in the degree of intermediate algebraic sets appearing during the process; this degree can be bounded by the Bézout bound recalled below, which is singly exponential in the number of variables. An implementation of this algorithm is available in the Kronecker library [96] of the Magma computer algebra system.

The complexity of algorithms for solving zero-dimensional systems usually depends on three important factors: the number of variables, the degrees of polynomials appearing during the computation and the number of solutions of the given system.

The highest degree appearing in the computation of Gröbner bases is known as the *degree of regularity*. Estimating this degree is an essential step in many complexity results for Gröbner

bases (see, e.g., [69, 7, 187, 65, 70]). The renowned Macaulay bound on the degree of regularity of a generic zero-dimensional system is due to Lazard [133].

Let  $f_1 = \dots = f_n = 0$  be a square zero-dimensional system in  $\mathbb{C}[x_1, \dots, x_n]$  of total degree  $\deg(f_i)$ . Under some mild assumptions, the degree of regularity of this system is bounded by

$$\sum_{i=1}^n (\deg(f_i) - 1) + 1.$$

On the other hand, Bézout bound provides an upper bound on the number of solutions. By Heintz's version of Bézout theorem [101, Theorem 1], the zero-dimensional system  $(f_1, \dots, f_s)$  given above has at most  $\deg(f_1) \dots \deg(f_s)$  complex solutions. Moreover, this bound is reached for a randomly generated dense system.

Positive-dimensional systems are systems with infinitely many solutions. These systems arise frequently in applications related to geometry or depending on parameters, for e.g., robotics [45, 37], computer vision [63] or geometry [204]. Developing methods for solving parametric polynomial systems is a challenging subject with many research directions (see, e.g., [134, 154, 76, 147]).

Over an algebraically closed field, an analogue of the zero-dimensional parametrization for algebraic sets of higher dimension is the *rational parametrization* introduced in [178].

Let  $\mathbf{f} \in \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  with variables  $\mathbf{x} = (x_1, \dots, x_n)$  and parameters  $\mathbf{y} = (y_1, \dots, y_t)$ . A rational parametrization of  $\mathbf{f}$  consists of

- A square-free polynomial  $w \in \mathbb{Q}(\mathbf{y})[u]$  where  $u$  is a new variable;
- A polynomial  $h \in \mathbb{Q}[\mathbf{y}]$ ;
- A sequence of polynomials  $(v_1, \dots, v_n) \subset \mathbb{Q}(\mathbf{y})[u]$  such that, for  $\eta \in \mathbb{C}^t$  and  $h(\eta) \neq 0$ ,  $\eta$  does not cancel any denominator in  $(w, v_1, \dots, v_n)$  and

$$V(\mathbf{f}(\eta, \cdot)) = \left\{ \left( \frac{v_1}{\partial w / \partial u}(\eta, \vartheta), \dots, \frac{v_n}{\partial w / \partial u}(\eta, \vartheta) \right) \mid w(\eta, \vartheta) = 0, \frac{\partial w}{\partial u}(\eta, \vartheta) \neq 0 \right\};$$

- A sequence  $(\lambda_1, \dots, \lambda_n) \in \mathbb{Q}^n$  such that

$$u \cdot w' = \sum_{i=1}^n \lambda_i \cdot v_i \text{ mod } w$$

The results of [178] provide a proof of existence of such a parametrization and an algorithm, called parametric geometric resolution, to compute it under some assumptions.

Intuitively, the rational parametrization provides a generic description for the solutions of  $\mathbf{f}(\eta, \cdot)$  when  $\eta$  ranges over  $\mathbb{C}^t$ . Note that one can choose the polynomial  $h \in \mathbb{Q}[\mathbf{y}]$  above in a way such that the number of complex solutions of  $\mathbf{f}(\eta, \cdot)$  is invariant as long as  $\eta$  does not cancel  $h$ . However, when considering real solutions, the behavior is more sophisticated.

**Example 1.1.2.** We consider the system given by

$$x_1^2 - x_2^2 - y^2 = 2x_2^2 + y^2 + y - 1 = 0,$$

where  $(x_1, x_2)$  are the variables and  $y$  is the parameter.

Let  $\Delta = (y^2 + y - 1)(y^2 - y + 1)$ . While  $\Delta \neq 0$ , this system always has exactly 4 distinct complex solutions:

$$\left( \pm \sqrt{\frac{y^2 - y + 1}{2}}, \pm \sqrt{\frac{-y^2 - y + 1}{2}} \right).$$

On the other hand, the number of distinct real solutions depends also on the sign of  $-y^2 - y + 1$ . It has 4 distinct real solutions when  $-y^2 - y + 1 > 0$  and no real solution when  $-y^2 - y + 1 < 0$ .

Therefore, to solve parametric polynomial systems over the reals, it requires more algorithmic results. This field is an important and active domain in both computer algebra and real algebraic geometry. In this thesis, we focus on two problems in this topic:

- One block quantifier elimination for systems of polynomial equations;
- Real root classification for parametric systems of polynomial equations.

Another problem which comes to our interest is to compute the *isolated points* of a semi-algebraic set. Such problems arise frequently in studying rigid mechanisms in material design [95, 75, 22, 23], which are modeled with semi-algebraic constraints. Naturally, isolated points of the semi-algebraic set under consideration are related to the rigidity. The mathematical foundations and potential applications of rigid materials are active research fields. In this thesis, we aim to solve specifically the following algorithmic problem:

- Computing the isolated points of a *real algebraic set*.

In what follows, we give more precise statements for these problems.

**One block quantifier elimination.** The most natural question in solving parametric polynomial systems is to ask for which parameter values the given system has at least one real solution.

From a geometric point of view, one can look at the vanishing set of the system in the space of all variables and parameters. Then the problem boils down to compute a description of the projection of this vanishing set on the parameter space.

By Tarski's theorem [192, Theorem 31], the projection of a given semi-algebraic set is also semi-algebraic. The *one block quantifier elimination* problem aims to compute a description by semi-algebraic formulas of this projection. This problem appears in various applications from a wide range of domains: program verification [120, 189], biology [157, 27, 38], economics [156], control theory [1, 2]. The precise statement is as follows.

Given a semi-algebraic formula  $\Psi(\mathbf{x}, \mathbf{y})$  in the variables  $\mathbf{x} = (x_1, \dots, x_n)$  and the parameters  $\mathbf{y} = (y_1, \dots, y_t)$ ,  $\Psi$  defines a semi-algebraic set  $\mathcal{S} \subset \mathbb{R}^{n+t}$ . Let  $\pi$  be the projection on the  $\mathbf{y}$ -coordinates, i.e.,

$$\pi : (\mathbf{x}, \mathbf{y}) \mapsto \mathbf{y}.$$

We arrive at computing a semi-algebraic formula  $\Phi(\mathbf{y})$  such that

$$\pi(\mathcal{S}) = \{\mathbf{y} \in \mathbb{R}^t \mid \Phi(\mathbf{y}) \text{ is true}\}.$$

In other words, we need to compute a quantifier-free semi-algebraic formula  $\Phi$  satisfying the equivalence below

$$\exists \mathbf{x} : \Psi(\mathbf{x}, \mathbf{y}) \iff \Phi(\mathbf{y}).$$

This formulation explains the name “one block quantifier elimination”.

In our work, we consider a weaker variant of the one block quantifier elimination problem for systems of polynomial equations as below.

Let  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$ . We denote by  $V(\mathbf{f}) \subset \mathbb{C}^{n+t}$  the set of complex solutions of

$$f_1 = \dots = f_s = 0.$$

We compute a semi-algebraic formula  $\Phi(\mathbf{y})$  satisfying

- $\{\mathbf{y} \in \mathbb{R}^t \mid \Phi(\mathbf{y}) \text{ is true}\} \subset \pi(V(\mathbf{f}) \cap \mathbb{R}^{n+t})$ ;
- The Lebesgue measure of  $\pi(V(\mathbf{f}) \cap \mathbb{R}^{n+t}) \setminus \{\mathbf{y} \in \mathbb{R}^t \mid \Phi(\mathbf{y}) \text{ is true}\}$  is zero in  $\mathbb{R}^t$ .

In practice, the parameters are usually given by approximate values containing numerical errors. Thus, this variant is usually sufficient for solving many applications.

**Example 1.1.3.** We consider for example the torus given by the following equation

$$(x^2 + y_1^2 + y_2^2 + 8)^2 = 36(y_1^2 + y_2^2).$$

Its projection on the  $(y_1, y_2)$ -coordinate can be written as a quantified formula

$$\exists x : (x^2 + y_1^2 + y_2^2 + 8)^2 - 36(y_1^2 + y_2^2) = 0,$$

which is then equivalent to the quantifier-free semi-algebraic formula

$$y_1^2 + y_2^2 \leq 16 \wedge y_1^2 + y_2^2 \geq 4.$$

For our variant, we can return the formula

$$y_1^2 + y_2^2 < 16 \wedge y_1^2 + y_2^2 > 4.$$

Note that one block quantifier elimination is a particular case of quantifier elimination in which one receives a formula with blocks of variables nested by quantifiers

$$\circ \mathbf{x}_1 \circ \mathbf{x}_2 \cdots \circ \mathbf{x}_\ell : \Omega(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell, \mathbf{y}),$$

where  $\mathbf{x}_1, \dots, \mathbf{x}_\ell, \mathbf{y}$  are blocks of variables and  $\circ$  is an interlaced sequence of universal and existential quantifiers  $\{\forall, \exists\}$ .

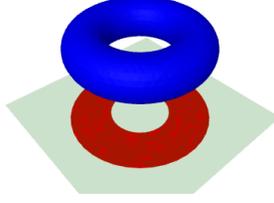


Figure 1.3: The projection of a torus is a semi-algebraic set.

**Real root classification.** A problem close to quantifier elimination is studied in [202, 203, 134, 143, 142] for the systems which have finitely many real solutions for almost every parameter values. For those systems, one can ask for which condition the system has a given number of real solutions.

As a consequence of Hardt's trivialization theorem [97, Sec. 4], those conditions on the parameters can be defined by semi-algebraic formulas. The computation of those semi-algebraic conditions is known as *real root classification* and appears in many applications such as robotics [134], computer vision [74, 64], physics [99], control theory [103, 119], etc.

Given  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  where  $\mathbf{x} = (x_1, \dots, x_n)$  are the variables and  $\mathbf{y} = (y_1, \dots, y_t)$  are the parameters, we denote by  $V(\mathbf{f}) \subset \mathbb{C}^{n+t}$  the set of complex solutions of

$$f_1 = \dots = f_s = 0$$

and  $\pi$  is the projection on the  $\mathbf{y}$ -space, i.e.,  $\pi : (\mathbf{x}, \mathbf{y}) \mapsto \mathbf{y}$ .

We assume that, for every parameter value  $\eta \in \mathbb{C}^t$  outside the vanishing set of some polynomial  $h \in \mathbb{C}[\mathbf{y}]$ , the system

$$f_1(\cdot, \eta) = \dots = f_s(\cdot, \eta) = 0$$

has finitely many complex solutions. By Hardt's triviality theorem [97, Sec. 4], there exists a finite collection of disjoint semi-algebraic sets  $\mathcal{S}_1, \dots, \mathcal{S}_\ell$  of  $\mathbb{R}^t$  such that

- The union  $\mathcal{S}_i$  is dense in  $\mathbb{R}^t$ ;
- The cardinality of  $\pi^{-1}(\eta) \cap V(\mathbf{f})$  is finite and invariant when  $\eta$  varies over each  $\mathcal{S}_i$ .

We aim to compute a list of tuples

$$\{(\Phi_1(\mathbf{y}), \eta_1, r_1), \dots, (\Phi_\ell(\mathbf{y}), \eta_\ell, r_\ell)\}$$

where

- $\Phi_i(\mathbf{y})$  are semi-algebraic formulas defining semi-algebraic sets  $\mathcal{S}_1, \dots, \mathcal{S}_\ell$  satisfying the above properties;

- $\eta_i$  is a point belong to  $\mathcal{S}_i$ ;
- $r_i$  is the number of real solutions of

$$f_1(\cdot, \eta_i) = \dots = f_s(\cdot, \eta_i) = 0.$$

A weaker output of the real root classification is the list of only the sample points and the number of real solutions:

$$\{(\eta_1, r_1), \dots, (\eta_\ell, r_\ell)\}.$$

This output is sufficient to identify which numbers of real solutions the input system possibly has.

**Example 1.1.4.** A complete real root classification of the torus in Example 1.1.3 is

$$\begin{aligned} 2 \text{ solutions} & : (y_1^2 + y_2^2 < 16) \wedge (y_1^2 + y_2^2 > 4), \\ 1 \text{ solution} & : (y_1^2 + y_2^2 = 16) \vee (y_1^2 + y_2^2 = 4), \\ 0 \text{ solution} & : (y_1^2 + y_2^2 > 16) \vee (y_1^2 + y_2^2 < 4). \end{aligned}$$

An admissible output of our variant can be

$$\begin{aligned} 2 \text{ solutions} & : (y_1^2 + y_2^2 < 16) \wedge (y_1^2 + y_2^2 > 4), \\ 0 \text{ solution} & : (y_1^2 + y_2^2 > 16) \vee (y_1^2 + y_2^2 < 4). \end{aligned}$$

**Computing isolated points of a real algebraic set.** Given  $f \in \mathbb{Q}[x_1, \dots, x_n]$ , the set of complex solutions of the equation  $f = 0$  is denoted by  $\mathcal{H}$ . A point  $\mathbf{x} \in \mathcal{H}$  is an isolated point of  $\mathcal{H} \cap \mathbb{R}^n$  if there exists  $r > 0$  such that

$$\mathcal{H} \cap \{\eta \in \mathbb{R}^t \mid \|\mathbf{x} - \eta\| < r\} = \{\mathbf{x}\},$$

where  $\|\cdot\|$  means the Euclidean norm in  $\mathbb{R}^n$ .

The set of isolated points of  $\mathcal{H} \cap \mathbb{R}^n$  is denoted by  $\mathcal{I}(\mathcal{H})$ . We aim to compute the following data, which allows to represent symbolically  $\mathcal{I}(\mathcal{H})$ ,

- A zero-dimensional parametrization  $\mathcal{C} = (w(u), v_1(u), \dots, v_n(u)) \subset \mathbb{Q}[u]$  such that

$$\mathcal{I}(\mathcal{H}) \subset \left\{ \left( \frac{v_1(\eta)}{w'(\eta)}, \dots, \frac{v_n(\eta)}{w'(\eta)} \right) \mid \eta \in \mathbb{R} : w(\eta) = 0 \right\};$$

- A set of disjoint intervals  $I_1, \dots, I_{|\mathcal{I}(\mathcal{H})|}$  of rational endpoints such that each  $I_i$  contains exactly one real root  $\eta_i$  of  $w$  and

$$\mathcal{I}(\mathcal{H}) = \left\{ \left( \frac{v_1(\eta_i)}{w'(\eta_i)}, \dots, \frac{v_n(\eta_i)}{w'(\eta_i)} \right) \mid 1 \leq i \leq |\mathcal{I}(\mathcal{H})| \right\}.$$

Note that the set of real solutions of a system of polynomial equations

$$f_1 = \cdots = f_s = 0$$

where  $f_1, \dots, f_s \in \mathbb{R}[x_1, \dots, x_n]$  coincides with the real solutions of a single equation

$$f_1^2 + \cdots + f_s^2 = 0.$$

Thus, the problem statement above covers all real algebraic sets.

**Example 1.1.5.** We consider the algebraic curve  $\mathcal{H}$  in  $\mathbb{C}^2$  defined by the equation

$$x_1^2 = x_2^3 - x_2^2.$$

The set of real isolated points  $\mathcal{I}(\mathcal{H})$  of  $\mathcal{H}$  contains the unique point  $\{(0, 0)\}$ .

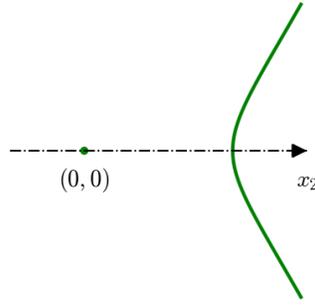


Figure 1.4: Real solution set of  $x_1^2 = x_2^3 - x_2^2$  with an isolated point  $(0, 0)$ .

As an example for our data representation, we take

$$\mathcal{C} = (u^2 - u, 0, u), \quad \mathcal{B} = \{[-1/2, 1/2]\}$$

to represent the set  $\mathcal{I}(\mathcal{H})$ . The zero-dimensional parametrization  $\mathcal{C}$  represents two points

$$(0, 0) = \left(0, \frac{u}{2u-1}\right)_{u=0} \quad \text{and} \quad (0, 1) = \left(0, \frac{u}{2u-1}\right)_{u=1}$$

in the given curve. The isolating box  $[-1/2, 1/2]$  means that the only root of  $u^2 - u = 0$  lying in  $[-1/2, 1/2]$  corresponds to the real isolated point of  $\mathcal{H}$ , which means the point  $(0, 0)$ .

## 1.2 State-of-the-art

In this section, we go through the prior works, which consist of the state-of-the-art on the complexity and software implementations for the aforementioned problems.

Throughout the thesis, we take as input elements from the field of rational numbers  $\mathbb{Q}$  which allows exact representation for symbolic computation. We measure only the arithmetic complexity of algorithms, i.e., the number of arithmetic operations  $+$ ,  $-$ ,  $\times$ ,  $\div$ , in the field  $\mathbb{Q}$ .

We use the standard Landau  $O$  notation for the complexity model:

- Let  $f : \mathbb{R}_+^\ell \mapsto \mathbb{R}_+$  be a positive function. We let  $O(f)$  denote the class of functions  $g : \mathbb{R}_+^\ell \rightarrow \mathbb{R}_+$  such that there exist  $C, K \in \mathbb{R}_+$  such that for all  $\|x\| \geq K$ ,  $g(x) \leq Cf(x)$ , where  $\|\cdot\|$  is a norm of  $\mathbb{R}^\ell$ .
- The notation  $O^\sim$  denotes the class of functions  $g : \mathbb{R}_+^\ell \rightarrow \mathbb{R}_+$  such that  $g \in O(f \log^\kappa(f))$  for some  $\kappa > 0$ .

### 1.2.1 Cylindrical algebraic decomposition

We start this related work section by recalling the cylindrical algebraic decomposition (CAD) and the algorithms that compute it. They are the first implemented tools that allow one to compute on semi-algebraic sets and still being used in computer algebra. These implementations are available in many computer algebra systems such as Maple/Mathematica or dedicated software like QEPCAD [25], REDLOG [183] or SYNRAC [72].

A cylindrical algebraic decomposition adapted to a given semi-algebraic set  $\mathcal{S} \subset \mathbb{R}^n$  is a partition of  $\mathcal{S}$  into connected cells which are homeomorphic to  $]0, 1[^i$  for some  $0 \leq i \leq n$ . These cells are stacked into a cylindrical structure by projections.

Since such a decomposition is quite exhaustive, it provides a lot of information and allows us to answer many algorithmic questions on semi-algebraic sets.

The first algorithm for computing CAD is due to Collins [41]. Since then, there have been enormous contributions to improve CAD, for which we can name the works in [148, 107, 149, 24] that improve the projection operator or the partial CAD [42], a variant of CAD that removes some redundant computations at each step.

Note that the CAD eliminates a single variable at each step and, after each elimination step, the degrees of involving polynomials grow quadratically. Therefore, on an input defined by  $s$  polynomials in  $n$  variables of degree at most  $D$ , computing CAD requires at most

$$(sD)^{2^{O(n)}}$$

arithmetic operations of the underlying field (see, for e.g., [50, 26]). In practice, this complexity is reached on randomly generated dense systems. Hence, the use of CAD algorithms is usually limited to only 4 variables for non-trivial problems.

In what follows, we discuss related works for each problem we consider.

### 1.2.2 One block quantifier elimination

Let  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  be our input and  $D$  be a bound on the degree of each  $f_i$ .

Historically, the works of Tarski and Seidenberg [192, 182] provide the first algorithm for solving quantifier elimination. However, this algorithm is not elementary recursive. The CAD algorithm of Collins [41] is the first implemented algorithm for solving quantifier elimination. Its arithmetic complexity for our problem is therefore

$$(sD)^{2^{O(n+t)}}.$$

In [197], Weispfenning proposed the use of his comprehensive Gröbner systems [198] for quantifier elimination. This approach was later developed in [190, 199, 54]. A comprehensive Gröbner system of a parametric polynomial system consists of a finite partition of parameter space and a set of Gröbner bases corresponding to each of those regions. Once a comprehensive Gröbner system is acquired, it provides a partition of the  $\mathbf{y}$ -space that allows one to apply a real root counting algorithm to each cell. However, the computation of comprehensive Gröbner systems is known to be expensive in practice and up to our knowledge, there is no complexity bound given for this computation.

Initiated in [93, 94], another class of quantifier elimination algorithms makes use of the block structure of quantifiers and the so-called critical point method. On a quantified formula of the form

$$\exists \mathbf{x}_\ell \forall \mathbf{x}_{\ell-1} \cdots \exists \mathbf{x}_1 : \Omega(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell, \mathbf{y})$$

with blocks of variables  $\mathbf{x}_1, \dots, \mathbf{x}_\ell, \mathbf{y}$ , these algorithms eliminate a whole block  $\mathbf{x}_i$  at every step, starting with the block

$$\exists \mathbf{x}_1 : \Omega(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell, \mathbf{y}),$$

through a reduction to dimension zero. More specifically, it computes a semi-algebraic formula  $\Phi(\mathbf{x}_1, \dots, \mathbf{x}_\ell, \mathbf{y})$  such that

- $\exists \mathbf{x}_1 : \Omega(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell, \mathbf{y})$  is logically equivalent to  $\exists \mathbf{x}_1 : \Phi(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell, \mathbf{y})$ .
- For every value  $\eta$  of  $(\mathbf{x}_2, \dots, \mathbf{x}_\ell, \mathbf{y})$ ,  $\Phi(\mathbf{x}_1, \eta)$  has finitely many real solutions in  $\mathbf{x}_1$ .

While the first property reduces the quantifier elimination on  $\Omega$  to  $\Phi$ , the first condition allows us to apply parametric real root counting algorithms to eliminate  $\mathbf{x}_1$  by considering  $(\mathbf{x}_2, \dots, \mathbf{x}_\ell, \mathbf{y})$  as parameters.

**Example 1.2.1.** *We take*

$$\Omega(x_1, x_2, y) = ((x_1^2 + x_2^2 + y^2 + 8)^2 - 36(x_2^2 + y^2) = 0)$$

*and solve the quantifier elimination problem*

$$\exists(x_1, x_2) : \Omega(x_1, x_2, y).$$

*The semi-algebraic set defined by*

$$\Phi(x_1, x_2, y) = ((x_1^2 + x_2^2 + y^2 + 8)^2 - 36(x_1^2 + x_2^2) = 0) \wedge (x_1 = 0)$$

satisfies the properties required above as can be seen in Fig. 1.5.

Geometrically,  $\Omega$  defines a torus  $\mathcal{V}$  in  $\mathbb{R}^3$  and  $\Phi$  defines a curve  $\mathcal{S}$  going around this torus. Let  $\pi : (x_1, x_2, y) \mapsto y$ . We have that

- Every fiber of  $\pi$  intersects the curve  $\mathcal{S}$  at finitely many points;
- The projections on  $y$  of  $\mathcal{V}$  and  $\mathcal{S}$  coincide.

Thus, solving the one block quantifier elimination on  $\Phi$  would also give an output for the same problem on  $\Omega$ .

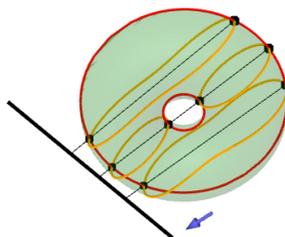


Figure 1.5: Reduction to real root classification.

Following the direction started by Grigor'ev and Vorobjov [93, 94], Heintz, Roy, Solernó [102], Renegar [164] and Basu, Pollack, Roy [8] introduced algorithms whose complexities are doubly exponential in only the number of blocks  $\ell$ . Particularly for one block of quantifier, the complexity of these algorithms is

$$s^{n+1} D^{O(nt)},$$

which is singly exponential in the number of variables (see [9, Algo. 14.6]). However, the techniques in use involves several infinitesimals which makes it challenging to be efficiently implemented. The only existing implementations for quantifier elimination are still those of CAD.

In spite of this tremendous progress, many important applications stay out of reach of the state-of-the-art software of quantifier elimination which are based on CAD. This motivates the studies of other variants, for instance, generic quantifier elimination [52, 183] and local quantifier elimination [53]. Generic quantifier elimination computes a quantifier-free formula that is equivalent to the input for almost all parameter values. On the other hand, local quantifier elimination returns an output which is equivalent to the input over a semi-algebraic set containing a target parameter values.

Recently, Hong and Safey El Din attempted to obtain a practical algorithm for a variant of quantifier elimination based on the critical point method in [109, 110]. Even though their algorithm applies under some assumptions, it shows an impressive performance in practice and solves various challenging stability analysis problems (6 indeterminates).

Motivated by these prior works, we aim to design practically efficient algorithms through a careful complexity driven approach. Moreover, we remark that the semi-algebraic formulas computed by the current algorithms for non-trivial problems are large and difficult to be evaluated. Therefore, we also want to have a compact representation of the output.

### 1.2.3 Real root classification

Again, we take  $(f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  as input where the degree of  $f_i$  is bounded by  $D$ .

Similar to quantifier elimination, a first approach for solving real root classification would be to compute a CAD of  $\mathbb{R}^{n+t}$  adapted to the input system. The cylindrical structure of the cells will imply that their projection on the parameters' space  $\mathbb{R}^t$  define semi-algebraic sets enjoying the properties required by the real root classification. However, the complexity

$$(sD)^{2^{O(n+t)}}$$

makes it difficult to use in practice.

The dedicated algorithms for real root classification were first designed for univariate polynomials with parametric coefficients [88, 143, 142]. These algorithms revisit different real root counting tools such as discriminant sequences or Sturm-Habicht sequences and apply them to parametric systems. To extend these results to multivariate systems, one relies on an algebraic elimination procedure to reduce to the univariate case (by, e.g., Gröbner bases or parametric geometric resolution). However, doing this usually introduces artificial singularities due to projections.

Another approach consists in computing a polynomial  $h \in \mathbb{Q}[\mathbf{y}]$  whose set of real solutions, denoted by  $V_{\mathbb{R}}(h)$ , contains the boundaries of semi-algebraic sets enjoying the properties required by the real root classification problem. More precisely, the number of real solutions of  $\mathbf{f}(\eta, \cdot)$  is invariant when  $\eta$  varies over each semi-algebraic connected component of  $\mathbb{R}^t \setminus V_{\mathbb{R}}(h)$ .

This is actually the direction followed by [203] and [134] in which such a polynomial  $h$  is called respectively *border polynomial* and *discriminant variety*. However, while [203] computes the *border polynomials* through triangular sets, [134] defines *discriminant variety* in a more geometric way and relies on eliminating procedures by Gröbner bases. Note that both [203, 134] do not discuss on the complexity aspect of their first steps.

Such a boundary immediately allows one to compute sample points for the weak output of real root classification using, e.g., [9, Chap. 13] whose complexity is  $\deg(h)^{O(t)}$  where  $\deg(h)$  is the total degree of  $h$ . On a randomly generated dense system  $\mathbf{f}$ , the polynomial  $h$  obtained using [134] defines the critical values of the restriction of  $\pi$  to  $V(\mathbf{f})$  whose degree is  $n(D-1)D^n$ . Hence, the step of computing sample points of  $\mathbb{R}^t \setminus V_{\mathbb{R}}(h)$  requires at most  $D^{O(nt)}$  arithmetic operations on these systems.

However, to obtain the full semi-algebraic descriptions, both [134] and [202] compute a CAD adapted to  $h \neq 0$ . Again, on the generic systems for which  $\deg(h) = n(D-1)D^n$ , the complexity for computing CAD of  $\mathbb{R}^t \setminus V_{\mathbb{R}}(h)$  lie in  $D^{n2^{O(t)}}$ .

An alternative method would be to use *parametric* roadmap algorithms to do such computations using, e.g., [9, Chap. 16] to compute semi-algebraic representations of the connected components of  $\mathbb{R}^t \setminus V(h)$ . Under the above extra assumptions, this would result in output formulas involving polynomials of degree bounded by  $(n(D-1)D^n)^{O(t^3)}$  using  $(n(D-1)D^n)^{O(t^4)}$  arithmetic operations (see [9, Theorem 16.13]). Note that the output degrees are by several orders of magnitude larger than  $n(D-1)D^n$  which bounds the degree of critical values of the restriction of  $\pi$  to  $\mathcal{V}$ .

### 1.2.4 Computing isolated points in a real algebraic set

Given as input a polynomial  $f \in \mathbb{Q}[x_1, \dots, x_n]$  of total degree  $D$ , we denote by  $\mathcal{H}$  the real algebraic set defined by  $f = 0$ .

As far as we know, there is no established algorithm dedicated to the computation of isolated points in a real algebraic set. However, effective real algebraic geometry provides subroutines from which such a computation could be done.

Note that the isolated points of  $\mathcal{H}$  coincide with the connected components of  $\mathcal{H}$  which are singletons. By Thom-Milnor bound [150, 193],  $\mathcal{H}$  has at most  $D(2D-1)^{n-1}$  connected components. This bound serves as an indicator for comparing the following approaches.

The first approach is to compute a CAD and to identify cells which correspond to isolated points using adjacency information (see e.g. [3]). The complexity of such a procedure is bounded by  $D^{2^{O(n)}}$ , which is doubly exponential in the number of variables  $n$ .

The algorithm for computing local dimension in [196] allows to compute isolated points in time  $D^{O(n^3)}$ .

A better approach is to formulate this problem by a quantified formula and use quantifier elimination algorithms to solve it. Using algorithms based on critical point method, for e.g. [9, Alg. 12.41], one obtains a complexity  $D^{O(n^2)}$ .

An alternative method is to use [9, Algorithm 12.16] to compute sample points in each connected component of  $\mathcal{H} \cap \mathbb{R}^n$  and then decide whether spheres, centered at these points, of infinitesimal radius, meet  $\mathcal{H} \cap \mathbb{R}^n$ . Note that these points are encoded with parametrizations of degree  $D^{O(n)}$  (their coordinates are evaluations of polynomials at the roots of a univariate polynomial with infinitesimal coefficients). Applying [9, Alg. 12.16] on this last real root decision problem would lead to a complexity  $D^{O(n^2)}$  since the input polynomials would have degree  $D^{O(n)}$ . One can also run [9, Alg. 12.16] modulo the algebraic extension used to define the sample points. That would lead to a complexity  $D^{O(n)}$  but this research direction requires modifications of [9, Alg. 12.16] since it assumes the input coefficients to lie in an *integral domain*, which is not satisfied in our case.

The topological nature of our problem is related to connectedness. Computing isolated points of  $\mathcal{H} \cap \mathbb{R}^n$  is equivalent to computing those connected components of  $\mathcal{H} \cap \mathbb{R}^n$  which are reduced to a single point. Hence, one considers computing *roadmaps*: these are algebraic curves contained in  $\mathcal{H}$  which have a non-empty and connected intersection with all connected components of the real set under study. Once such a roadmap is computed, it suffices to compute the isolated points

of a semi-algebraic curve in  $\mathbb{R}^n$ . This latter step is not trivial; as many of the algorithms computing roadmaps output either curve segments (see e.g., [11]) or algebraic curves (see e.g., [174]). Such curves are encoded through *rational parametrizations*, i.e., as the Zariski closure of the projection of the  $(x_1, \dots, x_n)$ -space of the solution set to

$$w(t, s) = 0, x_i = \frac{v_i(t, s)}{\partial w / \partial t}(t, s), \quad 1 \leq i \leq n$$

where  $w \in \mathbb{Q}[t, s]$  is square-free and monic in  $t$  and  $s$  and the  $v_i$ 's lie in  $\mathbb{Q}[t, s]$ . As far as we know, there is no published algorithm for computing isolated points from such an encoding.



Figure 1.6: A roadmap of a torus

Computing roadmaps started with Canny's (probabilistic) algorithm running in time  $D^{O(n^2)}$  on real algebraic sets. Later on, [173] introduced new types of connectivity results enabling more freedom in the design of roadmap algorithms. This led to [173, 11] for computing roadmaps in time  $(nD)^{O(n^{1.5})}$ . More recently, [10], still using these new types of connectivity results, provides a roadmap algorithm running in time  $D^{O(n \log^2 n)} n^{O(n \log^3 n)}$  for general real algebraic sets (at the cost of introducing a number of infinitesimals). This is improved in [174], for smooth bounded real algebraic sets, with a probabilistic algorithm running in time  $O((nD)^{12n \log_2 n})$ .

### 1.3 Contributions

This section provides an overview of our main contributions for the problems we consider. These contributions consist of new geometric results, the design of new algorithms with complexity improvements and implementations which outperform the software of the state-of-the-art. The main algorithmic tools for achieving these results are mainly based on the critical point method and the theory of Gröbner bases.

For the presented complexity results, we usually assume that our input polynomial sequence is generic. Informally, the genericity can be taken in the sense that the coefficients of input polynomials are considered as indeterminates that takes values in a suitable affine space; a rigorous definition is explained in Section 2.3.

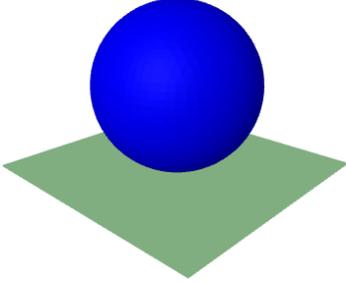
As our algorithm for the one block quantifier elimination problem makes use of the real root classification's one, we begin this section with the contributions for real root classification.

### 1.3.1 Real root classification

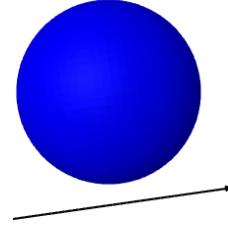
With M. Safey El Din, we design a new algorithm that solves the real root classification problem on an input  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  satisfying two assumptions:

- The polynomial sequence  $\mathbf{f}$  generates a radical ideal in  $\mathbb{Q}[\mathbf{x}, \mathbf{y}]$ ;
- For almost every  $\eta \in \mathbb{C}^t$ , the system  $f_1(\eta, \cdot) = \dots = f_s(\eta, \cdot) = 0$  has finitely many complex solutions.

**Example 1.3.1.** For example, the system  $\mathbf{f} = (x_1^2 + y_1^2 + y_2^2 - 1) \subset \mathbb{Q}[x_1, y_1, y_2]$  satisfies the both assumptions above. On the other hand, the system  $\mathbf{f} = (x_1^2 + x_2^2 + y_1^2 - 1) \subset \mathbb{Q}[x_1, x_2, y_1]$  fails to satisfy the second assumption.



(a)  $x_1^2 + y_1^2 + y_2^2 = 1$ .



(b)  $x_1^2 + x_2^2 + y_1^2 = 1$ .

We will see that this algorithm allows us to obtain a degree bound and an arithmetic cost which are better than the state-of-the-art for a wide range of inputs. To do that, we slightly generalize the notion of *Hermite quadratic forms*, a classical tool for counting solutions of zero-dimensional systems, to parametric systems. Originally introduced by Hermite [106] for counting real or complex solutions of univariate polynomials, Hermite quadratic forms were then extended in [160] to multivariate systems; their definition is as follows.

Given a zero-dimensional ideal  $I \subset \mathbb{Q}[\mathbf{x}]$ , Hermite quadratic forms operates on the finite dimensional  $\mathbb{Q}$ -vector space  $\mathbb{Q}[\mathbf{x}]/I$  by

$$\begin{aligned} \mathbb{Q}[\mathbf{x}]/I \times \mathbb{Q}[\mathbf{x}]/I &\rightarrow \mathbb{Q}, \\ (p, q) &\mapsto \text{trace}(\mathcal{L}_{p \cdot q}), \end{aligned}$$

where  $\mathcal{L}_{p \cdot q}$  denotes the endomorphism  $k \mapsto p \cdot q \cdot k$ .

Once a basis of the vector space  $\mathbb{Q}[\mathbf{x}]/I$  is fixed, such a quadratic form is represented by a symmetric matrix which is called a Hermite matrix. We will see in Subsection 4.4.3 how to obtain

a basis of  $\mathbb{Q}[\mathbf{x}]/I$  using a Gröbner basis of  $I$ , whose algorithmic properties depend also on the choice of the monomial ordering for the Gröbner basis.

The main property of Hermite matrices, given in [160, Theorem 2.1], states that: the number of distinct real (resp. complex) roots of the algebraic set defined by  $I$  equals the signature (resp. rank) of Hermite matrices.

**Example 1.3.2.** Given the ideal  $I = \langle x_1^2 + x_2x_1 + 2x_2 + 3, x_2^2 + 2x_1x_2 + 3x_1 + 1 \rangle$ , the Hermite matrices of  $I$  with respect to the bases  $B_1 = \{1, x_2, x_2^2, x_2^3\}$  and  $B_2 = \{1, x_2, x_1, x_2^2\}$  of  $\mathbb{C}[x_1, x_2]/I$  are respectively

$$H_1 = \begin{bmatrix} 4 & 5 & 97 & 818 \\ 5 & 97 & 818 & 7949 \\ 97 & 818 & 7949 & 74280 \\ 818 & 7947 & 74280 & 701998 \end{bmatrix} \quad \text{and} \quad H_2 = \begin{bmatrix} 4 & 5 & -1 & 97 \\ 5 & 97 & -49 & 818 \\ -1 & -49 & 27 & -338 \\ 97 & 818 & -338 & 7949 \end{bmatrix}.$$

The computation of the bases  $B_1, B_2$  and the Hermite matrices will be made clear in Subsection 4.4.3 when the preliminaries on Gröbner bases are fully introduced.

Since the both matrices have rank 4 and signature 2, we can deduce that the system

$$x_1^2 + x_2x_1 + 2x_2 + 3 = x_2^2 + 2x_1x_2 + 3x_1 + 1 = 0$$

has 4 distinct complex solutions and 2 distinct real solutions.

It is worth noting that the bit-sizes of coefficients in  $H_1$  is larger the ones in  $H_2$ , which makes the practical behaviors of these two matrices different in computations.

In our problem, we consider  $\mathbb{Q}(\mathbf{y})$  as the field of coefficients. The finiteness of generic fibers implies that  $\mathbf{f}$  generates a zero-dimensional ideal in  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$ . This allows us to carry out similar construction of Hermite's quadratic forms over  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$  to obtain what we call parametric Hermite quadratic forms. The matrices representing them are called parametric Hermite matrices.

**Example 1.3.3.** We consider the parametric system  $\mathbf{f} = (x_1^2 + x_2^2 - y_1, x_1x_2 + y_1x_2 + y_2)$ .

The parametric Hermite matrix associated to  $\mathbf{f}$  with respect to the basis  $B_1 = \{1, x_2, x_1, x_2^2\}$  of  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]/\langle \mathbf{f} \rangle$  is

$$\mathcal{H}_1 = \begin{bmatrix} 4 & 0 & -2y_1 & -2(y_1^2 - y_1) \\ * & -2(y_1^2 - y_1) & -4y_2 & -6y_1y_2 \\ * & * & 2(y_1^2 + y_1) & 2(y_1^3 - y_1^2) \\ * & * & * & 2(y_1^4 - 2y_1^3 + y_1^2 - 2y_2^2) \end{bmatrix}.$$

Whereas, using the basis  $B_2 = \{1, x_2, x_2^2, x_2^3\}$ , we obtain the parametric Hermite matrix

$$\mathcal{H}_2 = \begin{bmatrix} 4 & 0 & -2(y_1^2 - y_1) & -6y_1y_2 \\ * & * & * & 2(y_1^4 - 2y_1^3 + y_1^2 - 2y_2^2) \\ * & * & * & 10(y_1^3y_2 - y_1^2y_2) \\ * & * & * & -2(y_1^6 - 3y_1^5 + 3y_1^4 - 9y_1^2y_2^2 - y_1^3 + 3y_1y_2^2) \end{bmatrix}.$$

We also establish natural specialization properties for these parametric Hermite matrices. Hence, a parametric Hermite matrix, similar to its zero-dimensional counterpart, allows one to count respectively the number of roots at any parameters outside a strict algebraic sets of  $\mathbb{R}^t$  by evaluating the signature and rank of its specialization.

Based on this property, we derive from parametric Hermite matrices a polynomial  $w$  that plays the same role as the discriminant varieties [134] or the border polynomials [202]. Thus, one can compute the weak output for real root classification following the steps below.

- (a) Compute a parametric Hermite matrix  $\mathcal{H}$  associated to  $f \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  and a polynomial  $w_\infty \in \mathbb{Q}[\mathbf{y}]$  encoding the non-specialization locus of  $\mathcal{H}$ .

We rely on the theory of Gröbner bases to perform this step and also present some remarks to optimize the implementation of such a computation.

- (b) Compute a set of sample points  $\{\eta_1, \dots, \eta_\ell\}$  in the connected components of the semi-algebraic set of  $\mathbb{R}^t$  defined by  $w \neq 0$  where  $w$  is basically the product of  $\det(\mathcal{H})$  and  $w_\infty$ .

This is done through algorithms based on the critical point method (see e.g. [9, Chap. 12] and references therein) which are adapted to obtain practically fast algorithms following [171].

- (c) Compute the number  $r_i$  of real points in  $\mathcal{V} \cap \pi^{-1}(\eta_i)$  for  $1 \leq i \leq \ell$ .

This can be done by evaluating the signature of  $\mathcal{H}$  at the  $\eta_i$ 's.

To return semi-algebraic formulas, we follow a slightly different routine:

- (a) Compute a parametric Hermite matrix  $\mathcal{H}$  associated to  $f \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$ .
- (b) Compute a set of sample points  $\{\eta_1, \dots, \eta_\ell\}$  in the connected components of the semi-algebraic set of  $\mathbb{R}^t$  defined by  $\bigwedge_{i=1}^\delta M_i \neq 0$  where the  $M_i$ 's are the leading principal minors of  $\mathcal{H}$ . Again, this is done by the algorithm given in Section 5.2.
- (c) For  $1 \leq i \leq \ell$ , evaluate the sign pattern of  $(M_1, \dots, M_\delta)$  at the sample point  $\eta_i$ . From this sign pattern, we obtain a semi-algebraic formula representing the connected component corresponding to  $\eta_i$ .
- (d) Compute the number  $r_i$  of real points in  $\mathcal{V} \cap \pi^{-1}(\eta_i)$  for  $1 \leq i \leq \ell$ .

Note that the output formulas are encoded in determinantal forms and can be evaluated easily through the matrix  $\mathcal{H}$ .

**Example 1.3.4.** We continue with the system and the Hermite matrix  $\mathcal{H}_1$  in Example 1.3.3. The determinant of its parametric Hermite matrix is

$$w_{\mathcal{H}} = y_1^7 - 3y_1^6 - y_1^4 y_2^2 + 3y_1^5 + 20y_1^3 y_2^2 - y_1^4 + 8y_1^2 y_2^2 - 16y_2^4.$$

We notice that  $w_{\mathcal{H}}$  coincides exactly with the output returned by the procedure `DISCRIMINANTVARIETY` of Maple's commands `ROOTFINDING[PARAMETRIC]` that computes a discriminant variety [134].

Computing at least one point per connected component of the semi-algebraic set  $\mathbb{R}^3 \setminus V(w_{\mathcal{H}})$  using RAGlib gives us 12 points. We evaluate the signatures of  $\mathcal{H}$  specialized at those points and find that the input system can have 0, 2 or 4 distinct real solutions when the parameters vary.

Computing the leading principal minors of  $\mathcal{H}$ , we obtain

$$\begin{aligned} M_1 &= 4, \\ M_2 &= -8y_1(y_1 - 1), \\ M_3 &= -8(y_1^4 + y_1^3 - 2y_1^2 + 8y_2^2), \\ M_4 &= -16(y_1^7 - 3y_1^6 - y_1^4 y_2^2 + 3y_1^5 + 20y_1^3 y_2^2 - y_1^4 + 8y_1^2 y_2^2 - 16y_2^4). \end{aligned}$$

Since  $M_1$  is constant, we compute at least one point per connected component of the semi-algebraic set defined by

$$M_2 \neq 0 \wedge M_3 \neq 0 \wedge M_4 \neq 0.$$

The computation using RAGlib outputs a set of 22 sample points and finds the following realizable sign conditions of  $(M_2, M_3, M_4)$ :

$$[-1, -1, 1], [-1, -1, -1], [1, 1, 1], [1, -1, 1], [1, 1, -1], [1, -1, -1].$$

By evaluating the signature of  $\mathcal{H}$  at each of those sample points, we deduce the semi-algebraic formulas corresponding to every possible number of real solutions

$$\begin{aligned} 0 \text{ real root} &\rightarrow (M_2 > 0 \wedge M_3 < 0 \wedge M_4 > 0) \vee (M_2 < 0 \wedge M_3 < 0 \wedge M_4 > 0) \\ 2 \text{ real roots} &\rightarrow (M_2 > 0 \wedge M_3 < 0 \wedge M_4 < 0) \vee (M_2 < 0 \wedge M_3 < 0 \wedge M_4 < 0) \\ &\quad \vee (M_2 > 0 \wedge M_3 > 0 \wedge M_4 < 0) \\ 4 \text{ real roots} &\rightarrow (M_2 > 0 \wedge M_3 > 0 \wedge M_4 > 0). \end{aligned}$$

It can be seen in Example 1.3.3 that, when constructing the matrix  $\mathcal{H}$ , the degrees of polynomials involving in the parametric Hermite matrix depends on the choice of a monomial basis of the quotient ring of the ideal generated by  $\mathbf{f}$  in  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$ . This again depends on the monomial ordering used for Gröbner bases computations. In our work, we prioritize the so-called graded reverse lexicographical (grevlex) ordering (which acutally leads to the basis  $B_1$  in Example 1.3.3) whose interest for practical computations is explained in [13]. Notably, the following complexity statements of our algorithm are established for generic inputs using this grevlex ordering.

Given  $D \in \mathbb{N}$ , let  $\mathbf{f} = (f_1, \dots, f_n) \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  ( $s = n$ ) be a generic polynomial sequence such that  $\deg(f_i) \leq D$  and  $\mathfrak{D} = n(D - 1)D^n$ . We prove that

- i) Our algorithm computes the sample points for the output of real root classification within

$$O^{\sim} \left( \binom{t + \mathfrak{D}}{t} 8^t n^{2t+1} D^{2nt+n+2t+1} \right).$$

arithmetic operations in  $\mathbb{Q}$ .

- ii) When a complete output with semi-algebraic descriptions is required, this algorithm uses at most

$$O \sim \binom{t + \mathfrak{D}}{t} 8^t n^{2t+1} D^{3nt+2(n+t)+1}$$

arithmetic operations in  $\mathbb{Q}$  and the output involves polynomials of degree bounded by  $\mathfrak{D}$ .

We observe that the above degree bound is sharp for randomly generated dense systems.

The proof combines several techniques from the theory of Gröbner bases and algebraic geometry. First, it makes use of the genericity assumptions to establish some Noether position property for  $\mathbf{f}$ . This property allows us to connect the degree of entries in the parametric Hermite matrix with the degrees appearing in the grevlex Gröbner basis. Through the so-called Hilbert series of  $\mathbf{f}$ , we are able to control these latter degrees and deduce a bound for them.

Comparing with the state-of-the-art in [134, 202] which rely on computing a CAD with a complexity  $(sD)^{2^{O(t)}}$ , our algorithm has a complexity lying in  $D^{O(nt)}$ .

We implement our algorithm in MAPLE which uses the FGB library for Gröbner basis computation, MSOLVE for solving zero-dimensional systems and RAGLIB for computing sample points of semi-algebraic sets.

Table 1.8 below reports on the practical behavior of our algorithms. The columns HERMITE, CD and RRC represent respectively the timings of our implementation and the commands CellDecomposition (Discriminant variety approach) and RealRootClassification (Border polynomial approach) in MAPLE. The column DEG gives the highest degree of polynomials appearing in our output.

Our algorithm outperforms the other two state-of-the-art software and is able to solve various randomly generated dense systems and applications of real root classification which were not tractable. A particular application of our algorithms is the Kuramoto model [126]. This mathematical model, motivated by the behavior of chemical and biological systems, is used to describe the synchronization of coupled oscillators in many applications [185, 32, 18]. The Kuramoto models of 2 and 3 oscillators are studied carefully using computer algebra tools in [46]. In [99], the numerical solution for 4 oscillators is provided. Using our algorithm, we compute a similar answer as the one given by [99] which relies on numerical computation. Moreover, we obtain the semi-algebraic formulas for which the system has a given number of real solutions. As far as we known, this is the first symbolic solution for this application.

Especially, since the polynomials in our outputs are obtained as minors of parametric Hermite matrices, these matrices provide a compact determinantal representation of the output formulas, which then facilitates their evaluation. We illustrate this claim by reporting in Table 1.9 on the timings of these two different tasks for 1000 points  $\eta$ :

- Evaluating the signature of  $\mathcal{H}(\eta)$  (the column SIGN);
- Evaluating the principal minors of  $\mathcal{H}$  (the column MINORS);
- Solving specialized systems  $\mathbf{f}(\eta, \cdot)$  using MSOLVE, FGB and ROOTFINDING[ISOLATE] of Maple (the columns MSOLVE, FGB and ISOLATE).

System	$t$	$n$	$D$	HERMITE	DEG	CD	RRC
Dense	2	2	[2, 2]	.4 s.	8	.4 s.	1.1 s.
	2	2	[3, 2]	5 s.	18	1 m.	12 s.
	2	3	[2, 2, 2]	34 s.	24	17 m.	2 m.
	2	2	[3, 3]	3 m.	36	2 h.	4 m.
	3	2	[2, 2]	27 s.	8	36 s.	12 m.
	3	2	[3, 2]	3 h.	18	86 h.	37 h.
	3	3	[2, 2, 2]	32 h.	24	> 240 h.	> 120 h.
	3	2	[4, 2]	90 h.	32	> 240 h.	> 240 h.
4	2	[2, 2]	8 m.	8	> 240 h.	> 240 h.	
Kuramoto	3	6	[2, ..., 2]	86 h.	48	> 240 h.	> 240 h.

Figure 1.8: Timings of real root classification algorithms.

System	$t$	$n$	$D$	SIGN	MINORS	MSOLVE	FGB	ISOLATE
Dense	2	2	[2, 2]	.5 s	.2 s	2 s	12 s	33 s
	2	3	[2, 2, 2]	2 s	4 s	5 s	15 s	110 s
	2	2	[3, 3]	3 s	6 s	4 s	12 s	65 s
	2	2	[5, 2]	7 s	18 s	5 s	14 s	55 s
	2	2	[4, 3]	10 s	30 s	6 s	15 s	80 s
Dense	3	2	[2, 2]	.8 s	.4 s	2 s	10 s	16 s
	3	3	[2, 2, 2]	6 s	30 s	5 s	15 s	80 s
	3	2	[3, 3]	9 s	90 s	4 s	12 s	65 s

Figure 1.9: Timings for evaluating the formulas.

We note that evaluating the signatures of specialized Hermite matrices is faster than evaluating the minors. On the other hand, solving a specialized system would depend strongly on the number of variables  $n$  while evaluating the signatures depends on the number of parameters  $t$ . In the above examples where  $n = 2$  and  $t = 3$ , solving the specialized systems is better. Even though, the only library for solving polynomial systems is faster than evaluating the signatures on these examples is `MSOLVE`, which is highly optimized in C.

In a collaboration with D. Manevich and D. Plaumann, we apply our algorithm to an application coming from the theory of real algebraic curves. We consider the computation of simple totally real hyperplane sections, in which one asks for a given algebraic curve of degree  $\delta$ , whether there exists a hyperplane with real coefficients that intersect the curve at  $\delta$  distinct real points.

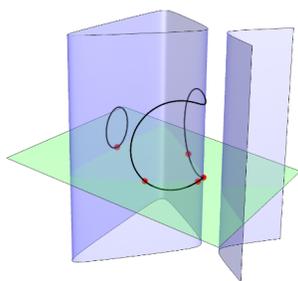
**Example 1.3.5.** For examples, let

$$\begin{aligned} f &= (x + 3)(x - y - 3)(x + y - 3) - 2, \\ g_1 &= x^2 + y^2 + z^2 - 10, \\ g_2 &= (x + 1)^2 + (y + 1)^2 + z^2 - 10. \end{aligned}$$

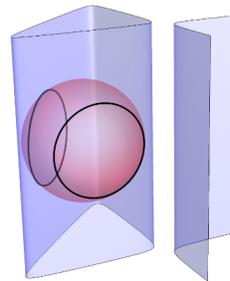
We take the algebraic curves defined by the intersection of the surface defined by  $f$  with respectively  $g_1$  and  $g_2$ . They are two real algebraic curves, each of which has two connected components. Our algorithm computes for the curve defined by  $f = g_1 = 0$  a simple totally real hyperplane section

$$x + \frac{43}{2000}y + \frac{131}{25}z + 9 = 0$$

while it concludes that the curve defined by  $f = g_2 = 0$  does not possess any simple totally real hyperplane section.



(a)  $f = g_1 = 0$



(b)  $f = g_2 = 0$

Figure 1.10: Examples for the (non)-existence of simple totally real hyperplane sections.

Such a question is motivated by the works on establishing explicit relations between quantities (degree, genus) of a given algebraic curve and the existence of totally real divisors [111, 152].

Taking the coefficients of the unknown hyperplane as parameters, this problem is naturally modeled as a real root classification problem. Using our real root classification algorithms, we illustrate a computational approach for deciding the existence of simple totally real hyperplane sections on some examples which then leads to the following findings:

1. There exist canonical curves  $X$  in  $\mathbb{P}^3$  with one or two ovals which do not allow simple totally real hyperplane sections.
2. There exists a curve  $X$  in  $\mathbb{P}^3$  of genus two and degree five having one oval which does not allow a simple totally real hyperplane section.
3. There are infinitely many plane quartics  $X$  with many ovals possessing a (complete) linear series of degree four which does not contain a totally real divisor.

### 1.3.2 One block quantifier elimination

This subsection presents new algorithmic and complexity results of a joint-work with M. Safey El Din for one block quantifier elimination problem.

Let  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  where  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{y} = (y_1, \dots, y_t)$  be the input. The set of complex solutions of

$$f_1 = \dots = f_s = 0$$

is denoted by  $\mathcal{V} \subset \mathbb{C}^n \times \mathbb{C}^t$  and  $\pi$  is the projection  $(\mathbf{x}, \mathbf{y}) \mapsto \mathbf{y}$ .

Our main result is a new probabilistic algorithm that takes an input  $\mathbf{f}$  that satisfies some regularity assumptions and computes a semi-algebraic formula  $\Phi$  defining a dense subset of the interior of  $\pi(\mathcal{V} \cap \mathbb{R}^{n+t})$ .

Note that, if almost every fiber  $\pi^{-1}(\eta) \cap \mathcal{V}$  is finite, any real root classification algorithm provides immediately a solution of one block quantifier elimination problem. With this in mind, we reduce the solving of one block quantifier elimination problem on  $\mathbf{f}$  to certain systems for which our real root classification can apply.

More precisely, we compute a polynomial system defining an algebraic set  $\mathcal{S} \subset \mathcal{V}$  such that there exists  $h \in \mathbb{Q}[\mathbf{y}]$  such that for  $\eta \in \mathbb{C}^t$  and  $h(\eta) \neq 0$ , the following holds:

- The fiber  $\pi^{-1}(\eta) \cap \mathcal{S}$  is finite.
- If  $\eta \in \mathbb{R}^t$ ,  $\pi^{-1}(\eta) \cap \mathcal{S} \cap \mathbb{R}^{n+t}$  is empty if and only if  $\pi^{-1}(\eta) \cap \mathcal{V} \cap \mathbb{R}^{n+t}$  is empty.

These two properties imply that  $\pi(\mathcal{S} \cap \mathbb{R}^{n+t})$  is a semi-algebraic subset of  $\pi(\mathcal{V} \cap \mathbb{R}^{n+t})$  such that  $\pi(\mathcal{V} \cap \mathbb{R}^{n+t}) \setminus \pi(\mathcal{S} \cap \mathbb{R}^{n+t})$  has zero Lebesgue measure in  $\mathbb{R}^t$ .

Hence, it remains to compute a quantifier-free semi-algebraic formula defining a dense subset in the interior of  $\pi(\mathcal{S})$ , for which we call to the real root classification algorithm introduced in the previous section.

Following this idea, our algorithm for one-block quantifier elimination problem proceeds through two main steps as follows.

- a) We compute a list of polynomial systems  $S_1, \dots, S_{d+1}$  in  $\mathbb{Q}[\mathbf{x}, \mathbf{y}]$  that satisfy
  - Each  $S_i$  generates a zero-dimensional ideal in  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$ ;
  - The union of the algebraic sets defined by  $S_i$  satisfies the property above.

This reduction is carried out by a parametric variant of the algorithm in [171] which actually computes at least one point per connected component of a regular real algebraic set. More specifically, this algorithm relies on the following geometric result [171, Theorem 2].

Let  $\pi_i$  be the projection  $(x_1, \dots, x_n) \mapsto (x_1, \dots, x_i)$ . For almost every  $A \in \text{GL}(n, \mathbb{Q})$ , the union of the sets

$$\text{crit}(\pi_i, \mathcal{V}^A) \cap \pi_{i-1}^{-1}((0, \dots, 0)), \quad 1 \leq i \leq d+1,$$

where  $\mathcal{V}^A = \{A^{-1} \cdot \mathbf{x} \mid \mathbf{x} \in \mathcal{V}\}$ , is finite and meets all connected components of  $\mathcal{V}^A \cap \mathbb{R}^n$ . The system  $S_i$  is taken as a defining system of  $\text{crit}(\pi_i, \mathcal{V}^A)$  using Jacobian criterion.

- b) For each system  $S_i$ , our real root classification algorithm outputs a semi-algebraic formula  $\Phi_i$  whose zero set is dense in the interior of the projection of real solutions of  $S_i$ .

Finally, we return

$$\Phi = \bigvee_{i=1}^{d+1} \Phi_i$$

as the final output of the one block quantifier elimination.

We illustrate our algorithm by the following example.

**Example 1.3.6.** We consider the polynomial  $f = x_1^2 + y_1x_2^2 + y_2x_2 + y_3$  in  $\mathbb{Q}[y_1, y_2, y_3][x_1, x_2]$ . Let  $\Delta = y_2^2 - 4y_1y_3$ . The projection of  $V(f) \cap \mathbb{R}^5$  on  $(y_1, y_2, y_3)$  is

$$(\Delta \geq 0 \wedge y_1 > 0) \vee (y_1 < 0) \vee (y_1 = 0 \wedge ((y_2 \neq 0) \vee (y_2 = 0 \wedge y_3 \leq 0))).$$

Applying the parametric variant of [171] for  $A$  taken as the  $3 \times 3$  identity matrix, we obtain 2 systems

$$W_1 = \{2y_1x_2 + y_2, f\} \quad \text{and} \quad W_2 = \{f, x_1\}.$$

For these systems, we compute  $w_{1,\infty} = w_{2,\infty} = y_1$  and the Hermite matrices:

$$H_1 = \begin{pmatrix} 2 & 0 \\ 0 & -2y_3 + y_2^2/(2y_1) \end{pmatrix}, \quad H_2 = \begin{pmatrix} 2 & -y_2/y_1 \\ -y_2/y_1 & (-2y_1y_3 + y_2^2)/y_1^2 \end{pmatrix}.$$

The sequences of leading principal minors are respectively  $[2, \Delta/y_1]$  and  $[2, \Delta/y_1^2]$ .

We compute then 4 points representing 4 connected components of the semi-algebraic set defined by  $y_1 \neq 0 \wedge \Delta \neq 0$ :

$$(1, 1/8, 0), (-1, 1/8, 0), (1, 1/8, 1/128), (-1, 1/8, -1/128).$$

The matrix  $H_2$  has non-zero signature over the first and second points, which both lead to the sign condition  $\Delta > 0 \wedge y_1^2 > 0$ . Thus, we have

$$\Phi_2 = (\Delta > 0 \wedge y_1^2 > 0) \wedge (y_1 \neq 0).$$

For  $H_1$ , non-zero signatures are satisfied at the first and fourth points. Evaluating the sign of  $\Delta$  and  $y_1$  at those points gives

$$\Phi_1 = ((\Delta > 0 \wedge y_1 > 0) \vee (\Delta < 0 \wedge y_1 < 0)) \wedge (y_1 \neq 0).$$

The final output is therefore  $\Phi = \Phi_1 \vee \Phi_2$ , which is equivalent to

$$\begin{aligned} \Phi &= (\Delta > 0 \wedge y_1 > 0) \vee (\Delta < 0 \wedge y_1 < 0) \vee (\Delta > 0 \wedge y_1 \neq 0) \\ &= (\Delta > 0 \wedge y_1 > 0) \vee (\Delta \neq 0 \wedge y_1 < 0). \end{aligned}$$

The difference between  $\pi(V(f) \cap \mathbb{R}^5)$  and the semi-algebraic set defined by  $\Phi$  is contained in  $(\Delta = 0) \vee (y_1 = 0)$  which is of zero measure in  $\mathbb{R}^3$ .

Controlling the complexity of this algorithm leads us to estimate the cost of classifying real solutions of the systems  $S_i$ . Similar to the complexity analysis of the real root classification algorithm, this requires a degree bound of the minors of the associated Hermite matrices to  $S_i$ .

Recall that the complexity result of our real root classification algorithm relies on the genericity of the input system. However, since the systems  $S_i$  are obtained as minors of some suitable Jacobian matrices of  $\mathbf{f}$ , they are not generic but equipped with a *determinantal structure*. Hence, we need to establish new complexity results for these structured systems.

When  $\mathbf{f}$  is generic, determinantal systems constructed from  $\mathbf{f}$  enjoy many nice properties (Cohen-Macaulay ring, explicit forms of Hilbert series, ...). In order to prove the complexity of our one block quantifier elimination algorithm, we take advantages of the properties of determinantal systems to extend the complexity proof of real root classification problem. Our complexity results are stated as follows.

Our algorithm for one block quantifier elimination, in case of success, computes a semi-algebraic formula  $\Phi$  defining a dense subset of the interior of  $\pi(\mathcal{V} \cap \mathbb{R}^{n+t})$ . The output formula involves polynomials of degree at most

$$\mathfrak{B} = D^s(D-1)^{n-s} \left( 2(n-s)(D-1) \binom{n-1}{s-2} + (n(D-2) + s) \binom{n-1}{s-1} \right).$$

The arithmetic cost of this algorithm is bounded by

$$O \sim \left( 8^t \mathfrak{B}^{3t+2} \binom{t+\mathfrak{B}}{t} \right).$$

Even though our complexity result lies in  $D^{O(nt)}$  as the ones based on critical point method (e.g., [9, Algo 14.6]), we obtain explicitly a degree bound on the output formulas and the constant in the big- $O$  notation in the exponent of the complexity. Especially, our degree bound  $\mathfrak{B}$  is observed to be sharp for generic inputs and, if  $s$  is fixed and  $D = 2$ ,  $\mathfrak{B}$  becomes polynomial in  $n$ .

We should emphasize that the other algorithms using critical point method are not implemented. The state-of-the-art software are based on the computation of CAD whose arithmetic complexity is  $(sD)^{2^{O(n+t)}}$ . Furthermore, our output formulas are evaluated easily through the parametric Hermite matrices.

Our algorithm is implemented in MAPLE, which calls to our real root classification algorithm's implementation. It uses the libraries FGB for algebraic elimination, MSOLVE for solving zero-dimensional system and RAGLIB for computing sample points.

Table 1.11 reports on the practical behavior of this implementation, comparing with quantifier elimination commands in MAPLE (QUANTIFIERELIMINATION) and MATHEMATICA (REDUCE) on randomly generated dense and sparse systems. It allows us to solve examples, both generic and non-generic, that are out of reach of these software (up to 8 indeterminates).

System	$t$	$n$	$s$	HERMITE	DEG	MAPLE	MATHEMATICA
Dense	2	3	2	4 s.	24	> 120 h.	> 120 h.
	2	4	2	1.5 m.	40	> 120 h.	> 120 h.
	2	5	2	20 m.	56	> 120 h.	> 120 h.
	2	6	2	3 h.	72	> 120 h.	> 120 h.
	2	7	2	8 h.	88	> 120 h.	> 120 h.
	3	3	2	1 m.	24	> 120 h.	> 120 h.
	3	4	2	20 m.	40	> 120 h.	> 120 h.
	3	5	2	7 h.	56	> 120 h.	> 120 h.
	3	6	2	24 h.	72	> 120 h.	> 120 h.
	4	3	2	30 m.	24	> 120 h.	> 120 h.
	4	4	2	46 h.	40	> 120 h.	> 120 h.
	5	3	2	14 h.	24	> 120 h.	> 120 h.
Sparse	3	3	2	40 s.	22	> 120 h.	> 120 h.
	3	4	2	15 m.	34	> 120 h.	> 120 h.
	3	5	2	15 m.	32	> 120 h.	> 120 h.
	4	3	2	20 m.	22	> 120 h.	> 120 h.
	4	4	2	20 m.	20	> 120 h.	> 120 h.

Figure 1.11: Timings of one block quantifier elimination algorithms.

### 1.3.3 Computing the isolated points of a real algebraic set

Given a polynomial  $f \in \mathbb{Q}[x_1, \dots, x_n]$  of total degree  $D$ , the set of complex solutions of  $f$  is denoted by  $\mathcal{H}$  and the set of isolated points of  $\mathcal{H} \cap \mathbb{R}^n$  is denoted by  $\mathcal{I}(\mathcal{H})$ .

With M. Safey El Din and T. de Wolff, we design several symbolic algorithms to compute the set  $\mathcal{I}(\mathcal{H})$ . To our knowledge, they are the first symbolic algorithms dedicated to this problem despite our restriction to the particular case of real algebraic sets.

All these algorithms share the first step of computing a finite set  $\mathcal{C}$  of points that contains  $\mathcal{I}(\mathcal{H})$  as a subset. We call these points the *candidates*. The set of candidates is encoded by a zero-dimensional parametrization

$$\mathcal{C} = (w, v_1, \dots, v_n)$$

where  $w, v_1, \dots, v_n \in \mathbb{Q}[u]$  and

$$\mathcal{C} = \left\{ \left( \frac{v_1(u)}{w'(u)}, \dots, \frac{v_n(u)}{w'(u)} \right) \mid w(u) = 0 \right\}.$$

This step can be done by computing at least one point per connected components of  $\mathcal{H}$  using critical point method algorithms. As we restrict to the case of single equation, we refer to the al-

gorithm [169] for this computation. The algorithm in [169] is based on the deformation technique which is formulated mathematically as follows.

Let  $\varepsilon$  be a transcendental element of  $\mathbb{R}$  such that  $\varepsilon < r$  for any positive  $r \in \mathbb{R}$ . Instead of considering directly the critical points of  $\mathcal{H}$ , one can compute the critical points of the smooth algebraic set  $\mathcal{H}_\varepsilon$  defined by  $f = \varepsilon$ . The critical points of  $\mathcal{H}$  can be obtained by taking the limits of the critical points of  $\mathcal{H}_\varepsilon$ . However, doing computation with infinitesimals is known to be inefficient in practice. Hence, in [169], an elimination procedure with Gröbner bases is used to avoid infinitesimals.

Once the zero-dimensional parametrization

$$\mathcal{C} = (w, v_1, \dots, v_n)$$

representing the candidates is computed, it remains to decide for each candidate whether it is a real isolated point of  $\mathcal{H}$ . We design in Chapter 8 several routines for this second step.

Our first approach is based on constructing an algebraic curve connecting the candidates that we specifically define. The identification of real isolated points is then reduced to decide the connectivity of some points on this curve. The construction of this curve uses the roadmap algorithm given in [174], which has the best known complexity bound  $(nD)^{O(n \log(n))}$  for constructing roadmaps. The complexity of this step leads to an arithmetic complexity bound of

$$(nD)^{O(n \log(n))}$$

for this approach.

The second approach proposed below leads to an arithmetic complexity  $D^{O(n)}$ . Moreover, we present several subroutines in our implementations to avoid as much as possible costly computations, in particular the computations with infinitesimals or over algebraic extensions.

To identify which candidates are actually isolated, a natural idea is to check whether a sphere of infinitesimal radius intersects the hypersurface  $\mathcal{H}$ . This approach has two drawbacks.

Recall that the set of candidates  $\mathcal{C}$  are encoded by a zero-dimensional parametrization

$$\mathcal{C} = (w, v_1, \dots, v_n),$$

in which the degree of the polynomials  $w, v_1, \dots, v_n$  is of order  $O(D^n)$ . This parametrization is taken as input to the decision problem. Therefore, using classical algorithms (for e.g. [9, Alg. 12.16]) for solving this problem leads to a complexity  $D^{O(n^2)}$ . Moreover, we also want to avoid the use of infinitesimals which appears in these algorithms for deformation.

A workaround to obtain a complexity  $D^{O(n)}$  is to carry out the computation over the quotient ring  $\mathbb{Q}[T]/\langle w(T) \rangle$ . This leads to solving polynomial systems over  $\mathbb{Q}[t]/\langle w(T) \rangle$  for which we rely on a variant of geometric resolution given in the appendix of [174]. It should be noted that extending the geometric resolution to this quotient ring is not obvious since the domain is only a product of fields (and not necessarily a field).

To remove the infinitesimals, we compute a sufficiently small rational radius  $e_0$  such that for each candidate  $\eta \in \mathfrak{C}$ ,  $\eta$  is an isolated point of  $\mathcal{H} \cap \mathbb{R}^n$  if and only if

$$\{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x} - \eta\| = e_0\} \cap \mathcal{H} = \emptyset.$$

The computation of such a value  $e_0$  can be done by computing the critical values of the distance function  $\|\bullet - \eta\|$ .

Once an appropriate value  $e_0$  is acquired, it remains to solve polynomial systems with rational coefficients. This is easily done by classical tools from semi-algebraic geometry.

Along with these main algorithms, we introduce several optimizations for avoiding as much as possible the expensive computations. The optimizations consist of two subroutines which are launched before the main routines of our algorithms.

- The first subroutine computes a set of points  $\mathfrak{C}_1$  such that  $\mathfrak{C}_1 \subset \mathcal{I}(\mathcal{H})$ .
- The second subroutine computes a set of points  $\mathfrak{C}_2$  such that  $\mathcal{I}(\mathcal{H}) \subset \mathfrak{C}_2 \subset \mathfrak{C}$ .

Moreover, when  $|\mathfrak{C}_1| = |\mathfrak{C}_2|$  or even  $|\mathfrak{C}_1| = |\mathfrak{C}|$ , they coincide with the set of real isolated points  $\mathcal{I}(\mathcal{H})$  and we obtain the output without running any heavy computation. This is actually the case for every example we consider.

Our algorithms are implemented in MAPLE. We use the FGB library for computing Gröbner bases to perform algebraic elimination required by our algorithms. Solving zero-dimensional systems for computing the candidates is done by MSOLVE and the subroutine for deciding emptiness of semi-algebraic sets calls to RAGLIB. We also used our C implementation for bivariate polynomial system solving (based on resultant computations) which we need to analyze connectivity queries in roadmaps.

We take sums of squares of  $n$  random dense quadrics in  $n$  variables (with a non-empty intersection over  $\mathbb{R}$ ); we obtain *dense quartics* defining a finite set of points. None of these examples could be solved by CAD algorithm in Maple within 10 days.

The implementations of our algorithms allow us to solve all of these examples. Timings for our algorithm using roadmaps are given in the column RM-ALGO below. The column APPROX-ALGO reports on the timings of our implementation of the algorithm using approximations. Note that this implementation takes into account the optimizations that we mention above.

Roadmaps are obtained as the union of critical loci of some maps in slices of the input variety [174]. We report on the highest degree of these critical loci in the column SRMP. The column SQRI reports on the maximum degree of the bivariate zero-dimensional system we need to study to analyze connectivity queries on the roadmap.

We also implemented [9, Alg. 12.16] using the FLINT C library with evaluation/interpolation techniques instead to tackle coefficients involving infinitesimals. This algorithm only computes sample points per connected components. *That implementation was not able to compute sample points of the input quartics for any of our examples.* We then report in the column [BPR] on the degree of the zero-dimensional system which is expected to be solved by [9, Alg. 12.16]. This is to be compared with the columns SRMP and SQRI.

$n$	RM-ALGO	SRMP	SQRI	APPROX-ALGO	[BPR]	MAPLE
4	50 s.	36	359	1 m.	7290	> 10 d.
5	8 h.	108	4 644	12 m.	65 610	> 10 d.
6	20 d.	308	47 952	7 h.	590 490	> 10 d.

Figure 1.12: Timings of computing real isolated points.

## 1.4 Organization of the thesis

This thesis is composed of two main parts.

Part I (Chapter 2 to Chapter 4) contains the preliminaries that will be used to present our contributions. Chapters 2 and 3 recall basic definitions and properties in algebraic geometry and the theory of Gröbner bases that will be used frequently throughout the thesis. Chapter 4 is dedicated to preliminaries on real algebraic geometry, especially the ones used for critical point method and real root counting. The results contained in this part are not original. There, we give the precise references for each result where a proof can be found.

Our contributions are presented in Part II of the thesis, which consists of five chapters below:

- Chapter 5 introduces our algorithms for solving the real root classification using parametric Hermite matrices.

The content of this chapter is also presented in the paper “[Solving parametric systems of polynomial equations over the reals through Hermite matrices](#)” (Huu Phuoc Le and Mohab Safey El Din) [136], which is under revision for Journal of Symbolic Computation.

- Chapter 6 contains our contributions for one-block quantifier elimination.

These results are also presented in the conference paper “[Faster One Block Quantifier Elimination for Regular Polynomial Systems of Equations](#)” (Huu Phuoc Le and Mohab Safey El Din) [137], published in the proceeding of ISSAC 2021, St. Petersburg, Russia.

- In Chapter 7, we show how to apply the algorithm presented in Chapter 5 to compute the totally real hyperplane sections on algebraic curves.

This work is extracted from the computational part of the paper “[Computing totally real hyperplane sections and linear series on algebraic curves](#)” (Huu Phuoc Le, Dimitri Manevich and Daniel Plaumann) [135] to be appeared in the journal *Le Matematiche*. These computations are my main contributions to this paper.

- Chapter 8 contains our works on the computation of the real isolated points of an algebraic hypersurface.

It consists of the results presented in the ISSAC 2020 paper “[Computing the real isolated points of an algebraic hypersurface](#)” (Huu Phuoc Le, Mohab Safey El Din and Timo de Wolff) [138] and also new works in an on-going collaboration with M. Safey El Din.

**Part I**

**Preliminaries**

# Chapter 2

## Basic notions of algebra and geometry

In this first chapter, we recall the preliminaries of commutative algebra and algebraic geometry which will be used in the next chapters. The notions and results presented in this chapter can be found with more details in the books by Eisenbud [56] and Cox, Little and O'Shea [48].

Throughout this chapter,  $R$  is a commutative ring whose zero and unit are denoted respectively by  $0_R$  and  $1_R$ .

### 2.1 Ideals

In classical algebraic geometry, we study the solutions of systems of polynomial equations over an algebraically closed field. These solution sets are related to the ideals of a polynomial ring. Therefore, we start this chapter by recalling the definition and some properties of ideals.

We start with the definition of ideals of a commutative ring  $R$ .

**Definition 2.1.1** (Ideals). *A subset  $I$  of  $R$  is called an ideal if and only if  $I$  is a subring of  $R$  and, for any  $r \in R$ ,  $rI \subset I$ .*

*The ideal generated by a subset  $S$  of  $R$  is*

$$\langle S \rangle = \{r_1 s_1 + \cdots + r_k s_k \mid r_1, \dots, r_k \in R, s_1, \dots, s_k \in S, k \in \mathbb{N}\}.$$

**Definition-Proposition 2.1.2.** *For an ideal  $I$  of the commutative ring  $R$ , we define the equivalence relation of elements of  $R$  below*

$$r \sim r' \quad \text{if and only if} \quad r - r' \in I.$$

*The set of equivalent classes  $[r]$  for  $r \in R$  is called the quotient of  $R$  by  $I$  and denoted as  $R/I$ .*

*It is equipped with a ring structure with the operations inherited from  $R$*

$$[r] + [r'] = [r + r'] \quad \text{and} \quad [r] \cdot [r'] = [r \cdot r'].$$

**Example 2.1.3.** *Let  $R = \mathbb{Z}$ , every ideal of  $R$  has the form  $\langle r \rangle$  for some  $r \in \mathbb{Z}$ . The quotient  $\mathbb{Z}/\langle r \rangle$  is a ring.*

*The quotient of the polynomial ring  $R[x, y]$  by the ideal  $\langle x, y \rangle$  is isomorphic to  $R$ .*

The following operations are defined for the ideals of  $R$ .

**Definition-Proposition 2.1.4** ([48, Chap. 4, Sec. 2, 3, 4]). *Let  $I$  and  $J$  be two ideals of  $R$ . Then the following subsets of  $R$  are also ideals of  $R$ :*

- *Sum:*  $I + J = \{f + g \mid f \in I, g \in J\}$ ;
- *Product:*  $IJ = \langle fg \mid f \in I, g \in J \rangle$ ;
- *Intersection:*  $I \cap J$ ;
- *Radical:*  $\sqrt{I} = \{f \mid \exists k \in \mathbb{Z}_+ \text{ such that } f^k \in I\}$ ;
- *Saturation:*  $I : J^\infty = \{f \in R \mid \exists k \in \mathbb{Z}_+ \text{ such that } fJ^k \subset I\}$ .

Note that  $IJ \subset I \cap J$  and the equality holds if  $I + J = R$ .

**Example 2.1.5.** Taking  $R = \mathbb{C}[x_1, x_2]$ , let  $I = \langle x_1 \rangle$  and  $J = \langle x_1x_2 \rangle$ . The ideal  $I \cap J = \langle x_1x_2 \rangle$  while  $IJ = \langle x_1^2x_2 \rangle$ .

Below are some frequently used definitions.

**Definition 2.1.6.** Let  $I$  be an ideal of  $R$ .

- $I$  is a maximal ideal if  $I \neq R$  and there is no ideal  $J$  of  $R$  such that  $I \subsetneq J \subsetneq R$ .
- $I$  is a prime ideal if  $I \neq R$  and whenever  $a, b \in R$  and  $ab \in I$ , then either  $a \in I$  or  $b \in I$ .
- $I$  is a primary ideal if for all  $f, g \in R$ ,  $fg \in I$  implies that  $f \in I$  or there is  $k \in \mathbb{Z}_+$  such that  $g^k \in I$ .
- $I$  is a radical ideal if  $I = \sqrt{I}$ .

**Definition 2.1.7** (Noetherian ring). A ring  $R$  is noetherian if and only if any increasing chain of ideals

$$I_0 \subseteq I_1 \subseteq \cdots \subseteq I_s \subseteq \cdots$$

is stationary.

**Example 2.1.8.** Any field is a noetherian ring as its only ideals are  $\langle 0 \rangle$  and  $\langle 1 \rangle$ .

Noetherian rings are also characterized by the following property.

**Proposition 2.1.9** ([56, Sec. 1.4.]). A ring  $R$  is noetherian if and only if every ideal of  $R$  is finitely generated, i.e., it admits a finite generating set.

The theorem below, known as Hilbert's basis theorem, is a key ingredient to ensure the termination of algorithms in commutative algebra.

**Theorem 2.1.10** (Hilbert's basis theorem, [56, Theorem 1.2]). If a ring  $R$  is noetherian then the polynomial ring  $R[x_1, \dots, x_n]$  is noetherian. As a consequence, every ideal of  $R[x_1, \dots, x_n]$  is finitely generated.

To each commutative ring  $R$ , one can associate the following notion of Krull dimension.

**Definition 2.1.11** (Krull dimension). *The height of a prime ideal  $\mathfrak{p}$ , denoted by  $\text{height}(\mathfrak{p})$ , is the supremum of all integers  $s$  such that there exists a chain of distinct prime ideals:*

$$\mathfrak{p}_0 \subsetneq \mathfrak{p}_1 \subsetneq \dots \subsetneq \mathfrak{p}_s = \mathfrak{p}.$$

*The Krull dimension of  $R$ , denoted by  $\dim R$ , is the supremum of the heights of all prime ideals of  $R$ .*

**Example 2.1.12.** *The Krull dimension of any field is 0 as its only prime ideal is  $\langle 0 \rangle$ .*

*The Krull dimension of  $R[x_1, \dots, x_n]$  is  $n$  as we have the maximal chain of prime ideals*

$$\langle 0 \rangle \subsetneq \langle x_1 \rangle \dots \subsetneq \langle x_1, \dots, x_n \rangle.$$

*The Krull dimension of the quotient ring  $R[x_1, x_2]/\langle x_1 - x_2^2 \rangle$  is 1 since  $R[x_1, x_2]/\langle x_1 - x_2^2 \rangle$  is isomorphic to  $R[x_2]$ .*

**Theorem 2.1.13** ([56, Chap. 8, Theorem A]). *Let  $R$  be an integral domain of finite Krull dimension. For any  $\mathfrak{p}$  be a prime ideal of  $R$ , we have*

$$\text{height}(\mathfrak{p}) + \dim R/\mathfrak{p} = \dim R$$

**Example 2.1.14.** *The ideal  $I = \langle x_1^2 - x_2, x_1 \rangle$  is a prime ideal of  $\mathbb{C}[x_1, x_2]$ . It has height 2 as we have the maximal ascending chain of prime ideals*

$$\langle 0 \rangle \subsetneq \langle x_1^2 - x_2 \rangle \subsetneq \langle x_1^2 - x_2, x_1 \rangle.$$

*On the other hand, we have that*

$$\mathbb{C}[x_1, x_2]/\langle x_1^2 - x_2, x_1 \rangle \simeq \mathbb{C}[x_1]/\langle x_1 \rangle \simeq \mathbb{C}$$

*Hence, its Krull dimension is 0. Thus, the equality from Theorem 2.1.13 holds:*

$$\text{height}(I) + \dim \mathbb{C}[x_1, x_2]/I = \dim \mathbb{C}[x_1, x_2].$$

The following definitions are essential in studying “local” properties of commutative rings.

**Definition 2.1.15** (Local ring). *A ring  $R$  is a local ring if it has a unique maximal ideal.*

**Example 2.1.16.** *Any field is a local ring and its unique maximal ideal is  $\langle 0 \rangle$ .*

The following construction leads to more interesting local rings.

**Definition 2.1.17** (Localization). *A subset  $S$  of  $R$  is called a multiplicative set if  $1_R \in S$  and for  $a, b \in S$ ,  $a \cdot b \in S$ . For a multiplicative set  $S$ , we can define an equivalence relation  $\sim$  on  $R \times S$  by  $(a, s) \sim (a', s')$  if and only if there is an element  $u \in S$  such that  $u(as' - a's) = 0$ .*

We denote the equivalence class of a pair  $(a, s) \in R \times S$  by  $\frac{a}{s}$ . The set of all equivalence classes

$$\{a/s \mid a \in R \text{ and } s \in S\}$$

is called the localization of  $R$  at  $R \setminus S$ , denoted by  $R_{R \setminus S}$ . This set  $R_{R \setminus S}$  is a ring with addition and multiplication given by

$$\frac{a}{s} + \frac{a'}{s'} = \frac{as' + a's}{ss'} \quad \text{and} \quad \frac{a}{s} \cdot \frac{a'}{s'} = \frac{aa'}{ss'}.$$

**Example 2.1.18.** For any prime ideal  $\mathfrak{p}$  of a commutative ring  $R$ , the set  $R \setminus \mathfrak{p}$  is a multiplicative set. The localization  $R_{\mathfrak{p}}$  of  $R$  at  $\mathfrak{p}$  is a local ring whose maximal ideal consists of all elements  $a/s$  with  $a \in \mathfrak{p}$  and  $s \in R \setminus \mathfrak{p}$ .

## 2.2 Affine algebraic sets

Let  $\mathbb{F}$  be a field and  $\mathbb{K}$  be an algebraically closed field containing  $\mathbb{F}$ . Affine algebraic sets of  $\mathbb{K}^n$  are defined as solution sets in  $\mathbb{K}^n$  of systems of polynomial equations of  $n$  variables with coefficients in  $\mathbb{K}$ .

**Definition 2.2.1.** Let  $S$  be a subset of  $\mathbb{F}[x_1, \dots, x_n]$  and  $\mathbb{L}$  be an extension of  $\mathbb{F}$ . The subset of  $\mathbb{L}^n$  at which the polynomials in  $S$  vanish, i.e.,

$$\{(x_1, \dots, x_n) \in \mathbb{L}^n \mid f(x_1, \dots, x_n) = 0 \text{ for any } f \in S\}$$

is called the affine algebraic set of  $S$  over  $\mathbb{L}$ , denoted by  $V_{\mathbb{L}}(S)$ .

Let  $\langle S \rangle$  be the ideal of  $\mathbb{L}[x_1, \dots, x_n]$  generated by  $S$ . It is easy to prove that  $V_{\mathbb{L}}(S) = V_{\mathbb{L}}(\langle S \rangle)$ .

Conversely, for an algebraic set  $V \subset \mathbb{K}^n$ , the subset of  $\mathbb{F}[x_1, \dots, x_n]$  of elements vanishing over  $V$ , i.e.,

$$\{f \in \mathbb{F}[x_1, \dots, x_n] \mid \text{for any } \eta \in V, f(\eta) = 0\},$$

is an ideal of  $\mathbb{F}[x_1, \dots, x_n]$  and is denoted as  $I(V)$ .

Recall that, by Hilbert's basis theorem (Theorem 2.1.10), the ideals of  $\mathbb{F}[x_1, \dots, x_n]$  are finitely generated.

The case where the ground field is an algebraically closed field is particularly important. In this case, we have the Hilbert Nullstellensatz theorem.

**Theorem 2.2.2** (Weak Nullstellensatz, [184, Prop. A.9]). *Let  $I$  be an ideal of  $\mathbb{K}[x_1, \dots, x_n]$ . The algebraic set  $V_{\mathbb{K}}(I)$  is empty if and only if  $I = \langle 1 \rangle$ .*

From an algorithmic point of view, deciding the emptiness of an ideal  $I$  can be done by testing whether 1 lies in  $I$ . We will see in the theory of Gröbner bases presented in Chapter 3 that this is a particular instance of the ideal membership problem which can be done effectively by computing *normal forms*.

The strong Nullstellensatz theorem below states that the radical  $\sqrt{I}$  is the set of polynomials that vanish over  $V_{\mathbb{K}}(I)$ . It is equivalent to the weak Nullstellensatz.

**Theorem 2.2.3** (Strong Nullstellensatz, [48, Chap. 4, Sec. 2, Theorem 6]). *Let  $I$  be an ideal of  $\mathbb{K}[x_1, \dots, x_n]$ . Then,*

$$I(V_{\mathbb{K}}(I)) = \sqrt{I}.$$

We note that the assumption of algebraic closedness of  $\mathbb{K}$  is crucial in Theorem 2.2.2. For instance,  $\langle x^2 + 1 \rangle$  is a proper ideal of  $\mathbb{R}[x]$  but  $V_{\mathbb{R}}(x^2 + 1) = \emptyset$ .

**Theorem 2.2.4** ([48, Chap. 4, Sec. 5, Theorem 11]). *Let  $\eta = (\eta_1, \dots, \eta_n) \in \mathbb{K}^n$ . Then we have that*

$$I(\{\eta\}) = \langle x_1 - \eta_1, \dots, x_n - \eta_n \rangle$$

*and this ideal is a maximal ideal of  $\mathbb{K}[x_1, \dots, x_n]$ .*

One can carry out the following operations on ideals of  $\mathbb{K}[x_1, \dots, x_n]$ .

**Proposition 2.2.5** ([48, Chap. 4, Sec. 3]). *Let  $I$  and  $J$  be two ideals of  $\mathbb{K}[x_1, \dots, x_n]$ . Then, we have*

- $V_{\mathbb{K}}(I + J) = V_{\mathbb{K}}(I) \cap V_{\mathbb{K}}(J)$ ,
- $V_{\mathbb{K}}(I \cap J) = V_{\mathbb{K}}(I \cdot J) = V_{\mathbb{K}}(I) \cup V_{\mathbb{K}}(J)$ ,
- $V_{\mathbb{K}}(I : J^\infty) = \overline{V_{\mathbb{K}}(I) \setminus V_{\mathbb{K}}(J)}$ .

As a corollary, we immediately have the proposition below.

**Proposition 2.2.6** ([48, Chap. 4, Sec. 3]). *The sets  $\emptyset$  and  $\mathbb{K}^n$  are algebraic sets of  $\mathbb{K}^n$ . The intersection and finite union of algebraic sets of  $\mathbb{K}^n$  are also algebraic sets of  $\mathbb{K}^n$ .*

Proposition 2.2.6 implies that the algebraic sets of  $\mathbb{K}^n$  form the closed sets of a topology of  $\mathbb{K}^n$ . This topology is called Zariski topology.

When  $\mathbb{K} = \mathbb{C}$ , this topology is much coarser than the strong (Euclidean) topology of  $\mathbb{C}^n$ . For example, the only Zariski closed proper subset of  $\mathbb{C}$  are finite sets of points.

The Zariski topology also has the following usual notions.

**Definition 2.2.7.** *Given any set  $S \subset \mathbb{K}^n$ , we denote by  $\overline{S}$  the Zariski closure of  $S$ , i.e., the smallest algebraic set of  $\mathbb{K}^n$  containing  $S$ . The set  $S$  is said to be Zariski dense in an algebraic set  $V$  if  $V$  coincides with the Zariski closure of  $S$ .*

*An algebraic set  $V \subset \mathbb{K}^n$  is irreducible if for any two algebraic sets  $V_1, V_2 \subset \mathbb{K}^n$  such that*

$$V = V_1 \cup V_2,$$

*then either  $V_1 = V$  or  $V_2 = V$ .*

The proposition below characterizes algebraically irreducible algebraic sets.

**Proposition 2.2.8** ([48, Chap. 4, Sec. 5, Prop. 3]). *If  $V$  is irreducible, the associated ideal  $I(V)$  is a prime ideal.*

Any algebraic set admits a unique decomposition into irreducible algebraic sets.

**Proposition 2.2.9** (Irreducible decomposition). *An algebraic set  $V$  has a unique decomposition (up to order) into irreducible algebraic sets*

$$V = \bigcup_{i=1}^s V_i$$

where  $V_i \not\subset V_j$  for any  $i \neq j$ .

The above decomposition of an algebraic set is translated to the *primary decomposition* of a defining ideal as follows.

**Proposition 2.2.10** (Primary decomposition, [48, Chap. 4, Sec. 8]). *Consider an ideal  $I \subset \mathbb{K}[x_1, \dots, x_n]$ . Then, there are primary ideals  $I_1, \dots, I_s \in \mathbb{K}[x_1, \dots, x_n]$  such that*

$$I = I_1 \cap \dots \cap I_s,$$

where  $\sqrt{I_i} \neq \sqrt{I_j}$  and  $\bigcap_{j \neq i} I_j \not\subset I_i$ .

We call the set  $(I_1, \dots, I_s)$  a *minimal primary decomposition* of  $I$ . Note that the ideal  $I$  can admit different minimal primary decompositions. However, those primary decompositions involve always the same number of primary ideals and the radicals  $\sqrt{I_1}, \dots, \sqrt{I_s}$ , which are prime ideals, are the same. These prime ideals are called the *associated primes* of  $\mathbb{K}[x_1, \dots, x_n]/I$ .

**Example 2.2.11.** *Given an ideal  $I = \langle x_1^4 - 2x_1^3 + x_1^2, x_1^2x_2 - 2x_1x_2 + x_2 \rangle$ , a minimal primary decomposition of  $I$  is*

$$I = \langle (x_1 - 1)^2 \rangle \cap \langle x_1^2, x_2 \rangle.$$

The algebraic set  $V_{\mathbb{C}}(I)$  is the union of the line  $x_1 = 1$  and the point  $(0, 0)$ . The radical  $\sqrt{I} = \langle x_1(x_1 - 1), x_2(x_1 - 1) \rangle$  and has the decomposition

$$\sqrt{I} = \langle x_1 - 1 \rangle \cap \langle x_1, x_2 \rangle.$$

The decomposition of  $V$  into irreducible components corresponds to the primary decomposition of the radical ideal  $I(V)$  associated to  $V$ :

$$I(V) = \bigcap_{i=1}^s I(V_i).$$

**Definition 2.2.12** (Coordinate ring). *Let  $V$  be an algebraic set in  $\mathbb{K}^n$  and  $I(V)$  is the radical ideal of  $\mathbb{K}[x_1, \dots, x_n]$  associated to  $V$ . The quotient ring  $\mathbb{K}[x_1, \dots, x_n]/I(V)$  is called the coordinate ring of  $V$ . We denote it by  $\mathbb{K}[V]$ .*

*Note also that  $\mathbb{K}[x_1, \dots, x_n]/I(V)$  is also equipped with a structure of  $\mathbb{K}$ -vector space.*

Intuitively, the coordinate ring  $\mathbb{K}[V]$  can be seen as the ring of functions from  $V$  to  $\mathbb{K}$  which coincide with a polynomial over  $V$ .

**Definition 2.2.13** (Dimension). *Let  $V$  be an irreducible variety. The dimension of  $V$  is defined as the Krull dimension of the coordinate ring  $\mathbb{K}[V]$  (see Definition 2.1.11).*

Let  $V$  be an algebraic set of dimension  $d$ . It is worth noting that the irreducible components of  $V$  can have different dimensions and the highest dimension among those irreducible components equals  $d$ . We have the following definition.

**Definition 2.2.14** (Local dimension). *The local dimension of a point  $\eta \in V$  is defined as the largest dimension among the irreducible components of  $V$  containing  $\eta$ .*

**Example 2.2.15.** *Let  $V \subset \mathbb{C}^3$  be the union of the plane  $V(x_1)$  and the line  $V(x_1, x_2)$ . The dimension at  $(0, 0, 0)$  is 2 while the dimension at any other point in  $V(x_1, x_2)$  is 1.*

**Definition 2.2.16.** *An algebraic set is equidimensional if it is the union of finitely many irreducible algebraic sets of the same dimension.*

The proposition below states every non-empty Zariski open subset of some affine space  $\mathbb{K}^n$  is dense in  $\mathbb{K}^n$ . It allows us to define the notion of *genericity* in the next section.

**Proposition 2.2.17** ([100, Chap. 1, Exercise 1.6]). *Let  $V$  be an irreducible algebraic set of dimension  $d$ . Then any proper Zariski closed subset of  $V$  has dimension at most  $d - 1$ . Therefore, any non-empty Zariski open subset of  $V$  is Zariski dense in  $V$ .*

## 2.3 Genericity and changes of variables

In this thesis, our algorithms rely on certain properties that hold for generic polynomial sequences. The definition of *genericity* is given below.

**Definition 2.3.1** (Genericity). *A property  $P$  over  $m$  free variables  $(u_1, \dots, u_m)$  which take values in  $\mathbb{K}$  is a boolean function*

$$\begin{aligned} P : \quad \mathbb{K}^m &\rightarrow \{true, false\} \\ (u_1, \dots, u_m) &\mapsto P(u_1, \dots, u_m). \end{aligned}$$

*We say that  $P$  is true generically if the subset of  $\mathbb{K}^m$  over which  $P(u_1, \dots, u_m)$  is true contains a non-empty Zariski open subset of  $\mathbb{K}^m$ .*

When  $\mathbb{K} = \mathbb{C}$ , since the set over which  $P$  does not hold is contained in a proper Zariski closed subset of  $\mathbb{C}^m$ , it has zero measure (with the usual Lebesgue measure of  $\mathbb{C}^m$ ).

For  $D \in \mathbb{N}$ , let  $\mathbb{K}[x_1, \dots, x_n]_{\leq D} \subset \mathbb{K}[x_1, \dots, x_n]$  be the set of polynomials of degree at most  $D$ . By considering the coefficients of a polynomial as coordinates,  $\mathbb{K}[x_1, \dots, x_n]_{\leq D}$  naturally has the structure of an affine space. In this thesis, our algorithms rely on certain properties that hold for *generic polynomial sequences*, i.e., generic properties over these spaces of polynomials.

**Example 2.3.2.** Given  $D_1, \dots, D_s \in \mathbb{N}$ , the following properties are generic for polynomial sequences  $\mathbf{f} = (f_1, \dots, f_s) \in \prod_{i=1}^s \mathbb{C}[x_1, \dots, x_n]_{\leq D_i}$ .

- The ideal  $\langle \mathbf{f} \rangle$  is radical.
- The algebraic set  $V_{\mathbb{C}}(\mathbf{f}) \subset \mathbb{C}^n$  is smooth.

These two statements can be proved by using Jacobian criterion (Theorem 2.4.4).

We will also see in Proposition 2.7.8 that, a class of polynomial sequences, called homogeneous regular sequences, is generic among all polynomial sequences.

Another use of genericity is through the changes of variables. In many algorithms, applying a random linear changes of variables to the input polynomial systems can ensure certain assumptions required by the algorithms in use (see, e.g., [87, 171, 178]). Since the changes of variables will be extensively used in this thesis, we introduce a dedicated notation for it.

We consider the polynomial ring  $\mathbb{F}[x_1, \dots, x_n]$ . Let  $\text{GL}(n, \mathbb{F})$  be the set of invertible matrices of size  $n \times n$  with entries in  $\mathbb{F}$ .

Let  $p \in \mathbb{F}[\mathbf{x}]$  be a polynomial. For any  $A \in \text{GL}(n, \mathbb{F})$ , we denote by  $p^A$  the polynomial  $p(A \cdot \mathbf{x}) \in \mathbb{F}[\mathbf{x}]$ . This notation applies also to a set of polynomials  $S \subset \mathbb{F}[\mathbf{x}]$

$$S^A = \{p^A \mid p \in S\}.$$

Given an algebraic set  $V \subset \mathbb{K}^n$ ,  $V^A$  denotes the algebraic set

$$V(\{p^A \mid p \in I(V)\}) = \{A^{-1} \cdot \mathbf{x} \mid \mathbf{x} \in V\}.$$

The set  $\text{GL}(n, \mathbb{K})$  is a non-empty Zariski open subset of the affine space of matrices of size  $n$ . A polynomial system  $\mathbf{f} \subset \mathbb{F}[x_1, \dots, x_n]$  satisfies a property  $P$  under a generic change of variables means that there exists a non-empty Zariski open subset  $\mathcal{A} \subset \text{GL}(n, \mathbb{K})$  such that for any  $A \in \mathcal{A}$ ,  $\mathbf{f}^A$  satisfies  $P$ .

**Example 2.3.3** ([14]). Let  $I$  be a zero-dimensional radical ideal of  $\mathbb{K}[x_1, \dots, x_n]$ . There exists a non-empty Zariski open subset  $\mathcal{A}$  of  $\text{GL}(n, \mathbb{K})$  such that for any  $A \in \mathcal{A}$ , the reduced Gröbner basis of  $I^A$  with respect to the lexicographic ordering  $x_1 \succ \dots \succ x_n$  has the form:

$$\{x_1 - g_1, \dots, x_{n-1} - g_{n-1}, \dots, g_n\},$$

where  $g_1, \dots, g_n$  lie in  $\mathbb{K}[x_n]$ .

In Chapter 5 and 6, we consider polynomials in two blocks of indeterminates  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{y} = (y_1, \dots, y_t)$ . In these cases, we might apply a linear change of variables that acts only on the variables  $\mathbf{x}$  and leave  $\mathbf{y}$  invariant. The matrices associated to these changes of variables form a subset denoted by  $\text{GL}(n, t, \mathbb{F})$  of  $\text{GL}(n + t, \mathbb{F})$ .

## 2.4 Tangent spaces and singularities

Similar to differential geometry, we attach to each point in an algebraic set a vector space which is called *tangent space*.

**Definition 2.4.1** (Tangent spaces). *Let  $V$  be an algebraic set of  $\mathbb{K}^n$  and  $(f_1, \dots, f_s)$  be a generating set of  $I(V)$ . The tangent space at a point  $\eta \in V$ , denoted by  $T_\eta(V)$  is the right kernel of the Jacobian matrix*

$$J = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_s}{\partial x_1} & \dots & \frac{\partial f_s}{\partial x_n} \end{bmatrix}$$

*evaluated at  $\eta$ .*

*This definition above does not depend on the choice of the generators  $f_1, \dots, f_s$ .*

One can define singular points through the rank of the Jacobian matrix, or equivalently, the codimension of tangent spaces.

**Definition 2.4.2** ([184, Chap. 2, Sec. 1.4, p. 94]). *Let  $V$  be an algebraic set and  $f_1, \dots, f_s \in \mathbb{K}[x_1, \dots, x_n]$  be a generating set of  $I(V)$ . We denote by  $J$  the Jacobian matrix of  $f_1, \dots, f_s$  with respect to  $x_1, \dots, x_n$ .*

*A point  $\eta \in V$  of local dimension  $d_\eta$  (see Definition 2.2.14) is a regular point of  $V$  if  $\text{rank } J = n - d_\eta$ . Otherwise, if  $\text{rank } J < n - d_\eta$ , then  $\eta$  is a singular point of  $V$ .*

In other words, one can say that a point  $\eta$  is singular if the local dimension of  $\eta$  is smaller than the dimension of  $T_\eta(V)$ . In fact, an algebraic set  $V$  cannot have too many singularities.

**Proposition 2.4.3** ([48, Chap. 9, Sec. 6, Theorem 8]). *The set of singularities of an algebraic set  $V$  is contained in a proper Zariski closed subset of  $V$ .*

The following theorem provides a tool to compute the singular points of a given algebraic set.

**Theorem 2.4.4** (Jacobian criterion, [56, Theorem 16.19]). *Let  $\mathbf{f} = (f_1, \dots, f_s)$  be a sequence of polynomials in  $\mathbb{K}[x_1, \dots, x_n]$ . Assume that at any point  $\eta$  of  $V(\mathbf{f})$ , the Jacobian matrix associated to  $\mathbf{f}$  has rank  $s$ . Then the ideal generated by  $\mathbf{f}$  is radical and the algebraic set  $V(\mathbf{f})$  is either empty or smooth and equidimensional of dimension  $n - s$ .*

**Example 2.4.5.** *We consider the prime ideal  $\langle x_1^3 - x_2^2 \rangle \subset \mathbb{C}[x_1, x_2]$ . By Theorem 2.1.13, we have that*

$$\dim \mathbb{C}[x_1, x_2] / \langle x_1^3 - x_2^2 \rangle = \dim \mathbb{C}[x_1, x_2] - \text{height}(\langle x_1^3 - x_2^2 \rangle) = 1.$$

*So, the algebraic set of  $\mathbb{C}^2$  defined by  $x_1^3 - x_2^2 = 0$  has dimension 1. The Jacobian matrix of  $x_1^3 - x_2^2$  with respect to  $x_1, x_2$  is written*

$$[3x_1^2 \quad -2x_2].$$

*The only point where the rank of this Jacobian matrix is smaller than 1 is  $(0, 0)$ . Thus,  $(0, 0)$  is the only singular point of  $V_{\mathbb{C}}(x_1^3 - x_2^2)$ .*

## 2.5 Morphisms between affine algebraic sets

Now we study morphisms between algebraic sets.

**Definition 2.5.1** (Polynomial morphism). *Let  $V \subset \mathbb{K}^n$  and  $W \subset \mathbb{K}^m$  be two algebraic sets. A map  $\varphi : V \rightarrow W$  is a polynomial morphism if there exist  $m$  polynomials  $\varphi_1, \dots, \varphi_m \in \mathbb{K}[x_1, \dots, x_n]$  such that*

$$\varphi(\eta) = (\varphi_1(\eta), \dots, \varphi_m(\eta))$$

for any  $\eta \in V$ .

Now one can define an equivalence between algebraic sets.

**Definition 2.5.2.** *Two algebraic sets  $V$  and  $W$  are isomorphic if there exists a polynomial morphism  $\varphi : V \rightarrow W$  such that  $\varphi$  is bijective and  $\varphi^{-1}$  is also a polynomial morphism.*

**Example 2.5.3.** *Let  $V = V(x_1 - x_2^2) \subset \mathbb{C}^2$ . The polynomial morphism  $\varphi : \mathbb{C} \rightarrow \mathbb{C}^2$ ,  $t \mapsto (t, t^2)$  is an isomorphism from  $\mathbb{C}$  to  $V$ , whose inverse morphism is the projection  $(x_1, x_2) \mapsto x_2$ .*

This is translated to the equivalence between coordinate rings.

**Theorem 2.5.4** ([48, Chap.4, Sec. 2, Theorem 7]). *Two algebraic sets  $V$  and  $W$  are isomorphic if and only if their coordinate rings  $\mathbb{K}[V]$  and  $\mathbb{K}[W]$  are isomorphic as rings.*

**Definition 2.5.5** (Dominant morphism). *The morphism  $\varphi$  is dominant if and only if the image of every irreducible component  $\mathcal{V}'$  of  $\mathcal{V}$  by  $\varphi$  is Zariski dense in  $\mathcal{W}$ , i.e.  $\overline{\varphi(\mathcal{V}')} = \mathcal{W}$ .*

**Example 2.5.6.** *The equation  $x_1x_2 = 1$  defines an algebraic curve in  $\mathbb{C}^2$ . The projection of this curve on the  $x_2$ -coordinate is  $\mathbb{C} \setminus \{0\}$ , whose Zariski closure is  $\mathbb{C}$ . Thus, this projection is a dominant morphism.*

The dimension of generic fibers of a dominant morphism is known by this theorem below.

**Theorem 2.5.7** (Fiber dimension theorem, [184, Theorem 1.25]). *Let  $\varphi$  be a dominant morphism from  $V$  to  $W$ . Then,*

- *For any point  $\eta \in W$ , the fiber  $\varphi^{-1}(\eta)$  has dimension at least  $\dim V - \dim W$ .*
- *There exists a Zariski open subset  $O$  of  $W$  such that  $O \subset \varphi(V)$  and, for any  $\eta \in O$ ,*

$$\dim \varphi^{-1}(\eta) = \dim V - \dim W.$$

Finally, we introduce the notion of critical points of a polynomial morphism. They correspond to the points, named critical values, of the arriving space whose fibers are singular algebraic sets. The precise definition is given below.

**Definition 2.5.8.** Let  $V \subset \mathbb{K}^n$  be an equidimensional algebraic set and  $\varphi : V \rightarrow W$  be a polynomial morphism.

A point  $\eta \in V$  is a critical point of the map  $\varphi$  if  $\eta$  is a regular point of  $V$  and the differential of  $\varphi$  at  $\eta$ ,  $d\varphi_\eta : T_\eta V \rightarrow T_{\varphi(\eta)} W$ , is surjective. The image by  $\varphi$  of a critical point is called a critical value.

The set of all critical points of the restriction of  $\varphi$  to  $V$  is denoted by  $\text{crit}(\varphi, V)$ .

One can compute the critical points of a smooth equidimensional algebraic set using the variant of Jacobian criterion below.

**Theorem 2.5.9** (Jacobian criterion, [174, Lemma A.2]). Let  $V \subset \mathbb{K}^n$  be an equidimensional algebraic set of dimension  $d$  and  $(f_1, \dots, f_s)$  be a generating set of the ideal  $I(V)$ .

Let  $\varphi$  be a polynomial morphism

$$\varphi : \begin{array}{ccc} \mathbb{K}^n & \rightarrow & \mathbb{K}^m, \\ \eta & \mapsto & (\varphi_1(\eta), \dots, \varphi_m(\eta)). \end{array}$$

A point  $\eta \in V$  is a critical point of  $\varphi$  if and only if  $\eta$  is a regular point of  $V$  and the Jacobian matrix associated to  $(f_1, \dots, f_s, \varphi_1, \dots, \varphi_m)$

$$\begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_s}{\partial x_1} & \cdots & \frac{\partial f_s}{\partial x_n} \\ \frac{\partial \varphi_1}{\partial x_1} & \cdots & \frac{\partial \varphi_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial \varphi_m}{\partial x_1} & \cdots & \frac{\partial \varphi_m}{\partial x_n} \end{bmatrix}$$

has rank less than  $n + m - d$  at  $\eta$ .

Note that when  $V$  is not equidimensional, the locus at which the Jacobian matrix has rank less than  $n + m - d$  contains all irreducible components of  $V$  of dimension less than  $d$ .

**Example 2.5.10.** We consider the algebraic curve  $V$  defined by

$$f = x_1^2 - x_2^3 + x_2^2.$$

and the projection  $\varphi : (x_1, x_2) \mapsto x_2$ . As the system

$$f = \frac{\partial f}{\partial x_1} = \frac{\partial f}{\partial x_2} = 0$$

has the unique solution  $(0, 0)$ , this point is the only singularities of  $V$ .

The Jacobian matrix of  $(f, \varphi)$  with respect to  $(x_1, x_2)$  is

$$J = \begin{bmatrix} 2x_1 & -3x_2^2 + 2x_2 \\ 0 & 1 \end{bmatrix}.$$

Requiring  $\text{rank } J < 2$  is equivalent to  $\det J = 2x_1 = 0$ . Thus, we obtain the zero-dimensional system

$$f = x_1 = 0$$

with two solutions  $(0, 0)$  and  $(0, 1)$ , among which  $(0, 1)$  is a regular point. Hence, we conclude that  $(0, 1)$  is a critical point of  $V$  with respect to  $\varphi$ .

The following algebraic version of Sard's theorem ensures that the set of critical values of a polynomial morphism is not Zariski dense in its target space.

**Theorem 2.5.11** (Sard's theorem, [174, Prop. B.2]). *Let  $V$  be an equidimensional algebraic subset of  $\mathbb{C}^n$  and  $\varphi : V \rightarrow \mathbb{K}^m$  be a polynomial mapping. Then, the set of critical values of  $\varphi$  is contained in a proper Zariski closed subset of  $\mathbb{K}^m$ .*

As the critical values are the images of the critical points, to compute them, one can eliminate some variables from the polynomial system defining the critical points obtained by Jacobian criterion (see Section 3.2). More precisely, let  $I$  be the ideal defining  $\text{crit}(\varphi, V)$  and  $z_1, \dots, z_m$  be new variables, the ideal

$$(I + \langle z_1 - \varphi_1, \dots, z_m - \varphi_m \rangle) \cap \mathbb{Q}[z_1, \dots, z_m]$$

defines the Zariski closure of the critical values of  $\varphi$ . Such computation can be done using the elimination theory of Gröbner bases recalled in Section 3.2.

We also use Thom's weak transversality theorem for proving certain properties of critical locus.

**Theorem 2.5.12** ([39, Theorem 3.7.4]). *Let  $\varphi : \mathbb{K}^n \times \mathbb{K}^t \rightarrow \mathbb{K}^m$  be a polynomial mapping. For  $\mathbf{y} \in \mathbb{K}^t$ , we define the partial application*

$$\begin{aligned} \varphi_{\mathbf{y}} : \mathbb{K}^n &\rightarrow \mathbb{K}^m, \\ \mathbf{x} &\mapsto \varphi(\mathbf{x}, \mathbf{y}). \end{aligned}$$

*Let  $\mathcal{X} \subset \mathbb{K}^n$  be a non-empty Zariski open subset such that  $\mathbf{0} \in \mathbb{K}^m$  is a regular value of  $\varphi$  restricted to  $\mathcal{X} \times \mathbb{K}^t$ . Then there exists a non-empty Zariski open subset  $\mathcal{Y} \subset \mathbb{K}^t$  such that for any  $\mathbf{y} \in \mathcal{Y}$ ,  $\mathbf{0} \in \mathbb{K}^m$  is a regular value of  $\varphi_{\mathbf{y}}$ .*

## 2.6 Projective algebraic sets

Let  $\mathbb{K}$  be an algebraically closed field. The projective space  $\mathbb{P}^n(\mathbb{K})$  is the set of equivalent classes of points in  $\mathbb{K}^{n+1} \setminus \{(0, \dots, 0)\}$ ,

$$[x_0 : x_1 : \dots : x_n] = \{(\lambda x_0, \dots, \lambda x_n) \mid \lambda \in \mathbb{K} \setminus \{0\}\}.$$

When the field  $\mathbb{K}$  is explicit from the context, we simply write  $\mathbb{P}^n$  for  $\mathbb{P}^n(\mathbb{K})$ .

A polynomial  $f \in \mathbb{K}[x_0, x_1, \dots, x_n]$  is homogeneous of degree  $D$  if its terms have the same degree  $D$ . Note that, if  $(x_0, \dots, x_n)$  is a solution of  $f$ , then every point  $(\lambda x_0, \dots, \lambda x_n)$  is also a solution of  $f$  in the affine space  $\mathbb{K}^{n+1}$ .

Therefore, we can take the solutions of homogeneous polynomials as points in the projective space  $\mathbb{P}^n$ . We have the following definition.

**Definition 2.6.1.** *A projective algebraic set is a subset of  $\mathbb{P}^n$  that is defined as a vanishing locus of a system of homogeneous polynomial equations.*

As for affine algebraic sets, we associate each projective algebraic set  $V \subset \mathbb{P}^n$  with an ideal

$$I(V) = \{f \in \mathbb{K}[x_0, x_1, \dots, x_n] \mid \text{for any } \eta \in V, f(\eta) = 0\}.$$

**Definition 2.6.2.** *An ideal  $I$  of  $\mathbb{K}[x_0, x_1, \dots, x_n]$  is called homogeneous if it can be generated by a set of homogeneous polynomials.*

**Proposition 2.6.3** ([48, Chap. 8, Sec. 2, Prop. 4]). *Let  $V \subset \mathbb{P}^n$  be a projective algebraic set. Then  $I(V)$  is a homogeneous ideal.*

The relation between projective algebraic sets and affine algebraic sets can be established through affine charts.

**Definition 2.6.4** (Affine charts). *The projective space can be decomposed into*

$$\mathbb{P}^n = \bigcup_{i=0}^n \mathbb{A}_i^n,$$

where

$$\mathbb{A}_i = \{[x_0, \dots, x_n] \in \mathbb{P}^n \mid x_i = 1\}.$$

The  $\mathbb{A}_i^n$ 's are called the affine charts of  $\mathbb{P}^n$ , each of which is an affine space  $\mathbb{K}^n$ .

The following constructions are used to relate projective spaces with their affine charts.

**Definition 2.6.5.** *We introduce the following notations:*

- For a polynomial  $p \in \mathbb{F}[x_1, \dots, x_n]$ , the homogenization  ${}^h p$  of  $p$  with respect to  $x_0$

$${}^h p = x_0^{\deg(p)} p\left(\frac{x_1}{x_0}, \dots, \frac{x_n}{x_0}\right).$$

- For a homogeneous polynomial  $q \in \mathbb{F}[x_0, x_1, \dots, x_n]$ , the dehomogenization  ${}^a q$  of  $q$  with respect to  $x_0$

$${}^a q = q(1, x_1, \dots, x_n).$$

Let  $I$  be an ideal of  $\mathbb{F}[x_1, \dots, x_n]$  and  $J$  be a homogeneous ideal of  $\mathbb{F}[x_0, x_1, \dots, x_n]$ .

- The homogenization ideal of  $I$  is the ideal

$${}^h I = \langle {}^h p \mid p \in I \rangle.$$

- The dehomogenization ideal of  $J$  is the ideal

$${}^a J = \langle {}^a q \mid q \in J \rangle.$$

These constructions lead to a relation between affine and projective algebraic sets.

Since the map

$$(x_1, \dots, x_n) \mapsto (1, x_1, \dots, x_n)$$

establishes a bijective correspondence between the affine chart  $\mathbb{A}_0$  of  $\mathbb{P}^n(\mathbb{K})$  and the affine space  $\mathbb{K}^n$ , we obtain easily an affine algebraic set from a projective one by dehomogenization as follows.

**Proposition 2.6.6** ([48, Chap.8, Sec.2, Exercise 9]). *Let  $W = V(f_1, \dots, f_s) \subset \mathbb{P}^n(\mathbb{K})$  be a projective algebraic set defined by homogeneous polynomials  $f_i \in \mathbb{K}[x_0, \dots, x_n]$ . Then, the affine algebraic set  $V({}^a f_1, \dots, {}^a f_s) \subset \mathbb{K}^n$  can be identified with the subset  $W \cap \mathbb{A}_0$  of  $W$ .*

From affine spaces to projective spaces, we go through the following definition of *projective closures*.

**Definition-Proposition 2.6.7** (Projective closure, [48, Chap. 8, Sec. 4, Prop. 7]). *Given an affine algebraic set  $V \subset \mathbb{K}^n$ , the projective closure of  $V$  is the smallest projective variety in  $\mathbb{P}^n(\mathbb{K})$  containing  $V$ .*

*It is also the projective algebraic set defined by the homogeneous ideal  ${}^h I(V)$  where  $I(V) \subset \mathbb{K}[x_1, \dots, x_n]$  is the ideal associated to  $V$ .*

**Proposition 2.6.8** ([48, Chap. 8, Sec. 4, Theorem 8]). *Let  $I$  be an ideal of  $\mathbb{K}[x_1, \dots, x_n]$ . Then the projective closure of  $V_{\mathbb{K}}(I)$  is the projective algebraic set associated to the homogenization  ${}^h I$  of  $I$ .*

Homogenizing a generating set  $(f_1, \dots, f_s)$  of an ideal  $I \subset \mathbb{K}[x_1, \dots, x_n]$  in general does not give the ideal  ${}^h I$ . We only have the inclusion

$$\langle {}^h f_1, \dots, {}^h f_s \rangle \subset {}^h I.$$

This can be illustrated in the example below.

**Example 2.6.9.** Let  $I = \langle f_1, f_2 \rangle = \langle x_2 - x_1^2, x_3 - x_1^3 \rangle$  and

$$J = \langle {}^h f_1, {}^h f_2 \rangle = \langle x_2 x_0 - x_1^2, x_3 x_0^2 - x_1^3 \rangle \subset \mathbb{R}[x_0, x_1, x_2, x_3].$$

Taking  $f_3 = f_2 - x_1 f_1 = x_3 - x_1 x_2 \in I$ , we obtain a homogeneous polynomial

$${}^h f_3 = x_0 x_3 - x_1 x_2 \in {}^h I.$$

Assume that  ${}^h f_3 \in J$ , then there exists  $A_1, A_2 \in \mathbb{R}[x_0, x_1, x_2, x_3]$  such that

$${}^h f_3 = A_1 \cdot {}^h f_1 + A_2 \cdot {}^h f_2.$$

Note that  ${}^h f_1, {}^h f_2$ , and  ${}^h f_3$  are homogeneous of degrees 2, 3, and 2 respectively. By looking at the homogeneous components of  $A_1$  and  $A_2$ , we deduce that  ${}^h f_3$  is a constant multiple of  ${}^h f_1$ . However, this is clearly false. Thus, by contradiction, we conclude that  ${}^h f_3 \notin J$  and  $J \neq {}^h I$ .

In Section 3.2, we will see that, by homogenizing the generating set

$$(x_1^2 - x_2, x_1 x_2 - x_3, -x_1 x_3 + x_2^2),$$

which is actually a Gröbner basis of  $I$ , we obtain a generating set of the projective closure of a given affine algebraic set.

## 2.7 Hilbert series and regular sequences

This section recalls the definition and some properties of Hilbert series and regular sequences. Hilbert series are generating series encoding many useful information of homogeneous ideals in polynomial rings, for instance, the Krull dimension of the given ideal, the degree of regularity, etc. On the other hand, ideals generated by regular sequences over a polynomial ring enjoy many nice properties. In particular, the Hilbert series of those ideals are explicitly known.

Therefore, these two notions of Hilbert series and regular sequences are used to analyze the complexities of many algorithms over polynomial rings, especially the algorithms for computing Gröbner bases (see, e.g., [69, 187, 65, 70]).

In Chapters 5 and 6, we will use Hilbert series to estimate the complexities of our algorithms.

Let  $\mathbb{F}$  be a field. We consider the decomposition of  $\mathbb{F}[x_0, \dots, x_n]$

$$\mathbb{F}[x_0, \dots, x_n] = \bigoplus_{D=0}^{\infty} \mathbb{F}[x_0, \dots, x_n]_D,$$

where  $\mathbb{F}[x_0, \dots, x_n]_D$  is a  $\mathbb{F}$ -vector space of homogeneous polynomials of total degree  $D$ . We call this decomposition the grading of  $\mathbb{F}[x_0, \dots, x_n]$  with respect to the total degree.

**Proposition 2.7.1** ([194, Sec. 1.2]). *Let  $I$  be a homogeneous ideal of  $\mathbb{F}[x_0, \dots, x_n]$ . Then the quotient ring  $\mathbb{F}[x_0, \dots, x_n]/I$  can be decomposed into*

$$\mathbb{F}[x_0, \dots, x_n]/I = \bigoplus_{D=0}^{\infty} \mathbb{F}[x_0, \dots, x_n]_D / (I \cap \mathbb{F}[x_0, \dots, x_n]_D),$$

where each  $\mathbb{F}[x_0, \dots, x_n]_D / (I \cap \mathbb{F}[x_0, \dots, x_n]_D)$  is a  $\mathbb{F}$ -vector space of finite dimension.

Using this grading of  $\mathbb{F}[x_0, \dots, x_n]$ , one defines the Hilbert series for the homogeneous ideal  $I$  as follows. This series provides many information of the associated homogeneous ideal.

**Definition 2.7.2** (Hilbert series). *Let  $I \subset \mathbb{F}[x_0, \dots, x_n]$  be a homogeneous ideal. The Hilbert series associated to  $I$  is defined as*

$$\text{HS}_I(z) = \sum_{D=0}^{\infty} \dim \mathbb{F}[x_0, \dots, x_n]_D / (I \cap \mathbb{F}[x_0, \dots, x_n]_D) \cdot z^D,$$

where the notion of dimension is taken for  $\mathbb{F}$ -vector spaces.

Let  $R$  be a commutative ring. A regular sequence over a ring  $R$  is defined as follows.

**Definition 2.7.3** (Regular sequence). *Let  $R$  be a commutative ring. A sequence of elements*

$$r_1, \dots, r_s \in R$$

is said to be a regular sequence if and only if for any  $0 \leq i \leq s$ ,  $r_{i+1}$  is not a zero-divisor of  $R / \langle r_1, \dots, r_i \rangle$ .

Using Definition 2.7.3 for the polynomial ring  $R = \mathbb{F}[x_0, \dots, x_n]$ , we define the notion of regular sequence for homogeneous polynomials.

**Definition 2.7.4** (Homogeneous regular sequence). *Let  $\mathbb{F}$  be a field. Given a homogeneous polynomial sequence  $(f_1, \dots, f_s) \subset \mathbb{F}[x_0, \dots, x_n]$  with  $s \leq n$ , we say that  $(f_1, \dots, f_s) \subset \mathbb{F}[x_0, \dots, x_n]$  is a regular sequence if for any  $1 \leq i \leq s$ ,  $f_i$  is not a zero-divisor in  $\mathbb{F}[x_0, \dots, x_n] / \langle f_1, \dots, f_{i-1} \rangle$ .*

Let  $\mathbb{K}$  be an algebraically closed field. Geometrically, the homogeneous regular sequences over  $\mathbb{K}[x_0, \dots, x_n]$  correspond to complete intersections defined as below.

**Definition 2.7.5** (Complete intersection). *Let  $V$  be a projective algebraic set of dimension  $d$  in projective space  $\mathbb{P}(\mathbb{K}^n)$ . We call  $V$  a complete intersection if the ideal of  $V$  can be generated by exactly  $n - d$  elements of  $\mathbb{K}[x_0, \dots, x_n]$ .*

Note that one can define affine regular sequences by simply removing the homogeneity assumption of  $(f_1, \dots, f_s)$  from the above definition. However, as explained in [6, Sec 1.7], many important properties that hold for homogeneous regular sequences are no longer valid for the affine ones using this definition. Therefore, we use [6, Definition 1.7.2] of affine regular sequences, which is more restrictive but allows us to preserve results similar to the homogeneous case.

**Definition 2.7.6** (Affine regular sequence). *Let  $\mathbb{F}$  be a field and  $(f_1, \dots, f_s) \subset \mathbb{F}[x_1, \dots, x_n]$  where  $s \leq n$ . We say that  $(f_1, \dots, f_s) \subset \mathbb{F}[x_1, \dots, x_n]$  is an affine regular sequence if the homogeneous parts of highest degree of the  $f_i$ 's form a homogeneous regular sequence.*

An important property of a regular sequence  $(f_1, \dots, f_s)$  is that the explicit form of the Hilbert series associated to  $(f_1, \dots, f_s)$  of  $I$  is known.

**Proposition 2.7.7** ([6, Prop. 1.7.4]). *Let  $I = \langle f_1, \dots, f_s \rangle$  be a homogeneous ideal of  $\mathbb{F}[x_0, \dots, x_n]$ . The Hilbert series of  $\mathbb{F}[x_0, \dots, x_n]/I$  satisfies the inequality (coefficient-wise)*

$$\text{HS}_I(z) \geq \frac{\prod_{i=1}^s (1 - z^{\deg(f_i)})}{(1 - z)^n}.$$

*The equality occurs if and only if  $(f_1, \dots, f_s)$  forms a homogeneous regular sequence.*

Now we consider polynomial ring  $\mathbb{K}[x_0, \dots, x_n]$  over an algebraically closed field  $\mathbb{K}$ . Since the set  $\mathbb{K}[x_0, \dots, x_n]_D$  of all homogeneous polynomials of total degree  $D$  over  $\mathbb{K}$  can be identified as an affine space  $\mathbb{K}^{\binom{D+n}{n}}$  (see Section 2.3). The proposition below states that the set of regular sequences are Zariski dense among all polynomial sequences.

**Proposition 2.7.8** ([159]). *Fixing  $D_1, \dots, D_s \in \mathbb{N}$ . The set of homogeneous regular sequences of  $\prod_{i=1}^s \mathbb{K}[x_0, \dots, x_n]_{D_i}$  contains a non-empty Zariski open subset of  $\prod_{i=1}^s \mathbb{K}[x_0, \dots, x_n]_{D_i}$ .*

## 2.8 Cohen-Macaulay rings and determinantal ideals

Given a homogeneous regular sequence  $(r_1, \dots, r_s) \subset \mathbb{F}[x_0, \dots, x_n]$ , the number of elements  $s$  is necessarily smaller than the number of variables  $n + 1$ . However, in practice, we encounter frequently over-determined polynomial systems where the number of equations is larger than the number of variables. Particularly, such systems appear in the computation of critical points using Jacobian criterion (see Section 2.5), where the equations are obtained from many minors of a Jacobian matrix. The properties of regular sequences are no longer applicable for estimating the complexity for computing Gröbner bases on these systems. Hence, we will need the following notion of Cohen-Macaulay rings to handle a wider class of polynomial sequences.

**Definition 2.8.1** (Cohen-Macaulay ring). *Let  $I$  be a proper ideal of a commutative Noetherian ring  $R$ . We define the depth of  $I$ , denoted by  $\text{depth}(I)$ , as the length of any maximal regular sequence in  $I$  considered as a ring on itself.*

*A ring  $R$  such that, for every maximal ideal  $\mathfrak{m}$  of  $R$ ,*

$$\text{depth}(\mathfrak{m}) = \text{height}(\mathfrak{m})$$

*is called a Cohen-Macaulay ring.*

**Example 2.8.2.** We consider the matrix  $U = (u_{k,\ell})_{1 \leq i \leq k, 1 \leq j \leq \ell}$  where the  $u_{i,j}$ 's are indeterminates. For any  $r \leq \min\{k, \ell\}$ , the set

$$\{\eta \in \mathbb{C}^{k \times \ell} \mid \text{rank } U(\eta) < r\}$$

is an algebraic set of  $\mathbb{C}^{k \times \ell}$  defined by the simultaneous vanishing of the  $r$ -minors of  $U$ . This set is called a determinantal variety and the ideal generated by the  $r$ -minors of  $U$  is called a determinantal ideal. The quotient ring of this determinantal ideal is a Cohen-Macaulay ring (see [28, Cor. 2.8]).

The property that determinantal ideals are Cohen-Macaulay is used for proving complexity results for computing critical points of a polynomial function in [69, 187].

The following proposition gives another equivalent characterization for a Cohen-Macaulay ring through its localizations.

**Proposition 2.8.3** ([56, Prop. 18.8]). *A ring  $R$  is Cohen-Macaulay if and only if, for every prime ideal  $\mathfrak{p}$  of  $R$ , the localization  $R_{\mathfrak{p}}$  of  $R$  at  $\mathfrak{p}$  is Cohen-Macaulay.*

**Example 2.8.4.** *The quotient ring  $R = \mathbb{K}[x_1, x_2] / \langle x_1^2, x_1x_2 \rangle$  is not a Cohen-Macaulay ring since the localization  $R_{\mathfrak{p}}$  of  $R$  at the maximal ideal  $\mathfrak{p} = \langle x_1, x_2 \rangle$  is not Cohen-Macaulay.*

*Indeed, let  $\mathfrak{m}$  denote the unique maximal ideal  $\langle x_1, x_2 \rangle$  of  $R_{\mathfrak{p}}$ . Every  $f \in \mathfrak{m}$  is a zero-divisor since  $x_1 \neq 0$  and  $x_1f = 0$  (as  $x_1^2 = x_1x_2 = 0$ ). So,  $\mathfrak{m}$  does not contain any regular sequence, and therefore,  $\text{depth}(\mathfrak{m}) = 0$ . Thus,  $R_{\mathfrak{p}}$  is not Cohen-Macaulay.*

Cohen-Macaulay rings enjoy many nice properties.

**Proposition 2.8.5** ([56, Prop. 18.9]). *A ring  $R$  is Cohen-Macaulay if and only if  $R[x_1, \dots, x_n]$  is Cohen-Macaulay.*

**Proposition 2.8.6** ([56, Prop. 18.13]). *Let  $R$  be a Cohen-Macaulay ring. If  $I = \langle r_1, \dots, r_s \rangle$  is an ideal generated by  $s$  elements in  $R$  such that  $\text{height}(I) = s$  (the largest value), then  $R/I$  is a Cohen-Macaulay ring.*

**Theorem 2.8.7** (Unmixedness theorem, [56, Cor. 18.14]). *Let  $R$  be a ring. If  $I = \langle r_1, \dots, r_s \rangle$  is an ideal generated by  $s$  elements such that  $\text{height}(I) = s$ , then all minimal primes of  $I$  have codimension  $s$ . If  $R$  is Cohen-Macaulay, then every associated prime of  $I$  is minimal over  $I$ .*

## 2.9 Noether position and properness

**Definition 2.9.1** (Integral extension). *Let  $S$  be a ring containing  $R$ . An element  $u \in S$  is integral over  $R$  if there exist  $d \in \mathbb{N}$  and  $r_0, \dots, r_{d-1} \in R$  such that*

$$u^d + u_{d-1}u^{d-1} + \dots + u_0 = 0.$$

*If all elements of  $S$  are integral over  $R$ ,  $S$  is called an integral extension of  $R$ .*

**Example 2.9.2.** The variable  $x_1$  is integral over the quotient ring  $\mathbb{C}[x_1, x_2]/\langle x_1^2 - x_2 - 1 \rangle$  as the monic polynomial  $x_1^2 - x_2 - 1$  is 0 in this ring.

**Definition 2.9.3** (Noether position). Let  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{F}[x_1, \dots, x_{s+t}]$ . The variables  $(x_1, \dots, x_s)$  are in Noether position with respect to the ideal generated by  $\mathbf{f}$  if the following properties are satisfied:

- $\mathbb{F}[x_{s+1}, \dots, x_{s+t}] \cap \langle \mathbf{f} \rangle = \langle 0 \rangle$ , which implies  $\mathbb{F}[x_{s+1}, \dots, x_{s+t}] \subset \mathbb{F}[x_1, \dots, x_{s+t}]/\langle \mathbf{f} \rangle$ .
- The canonical images of  $x_1, \dots, x_s$  in the quotient algebra  $\mathbb{F}[x_1, \dots, x_{s+t}]/\langle \mathbf{f} \rangle$  are integral over  $\mathbb{F}[x_{s+1}, \dots, x_{s+t}]$ .

For homogeneous ideals, Noether position is strongly related to the regularity.

**Proposition 2.9.4** ([194, Prop. 1.44]). Let  $\mathbf{f} = (f_1, \dots, f_s)$  be a sequence of homogeneous polynomials in  $\mathbb{F}[x_1, \dots, x_{s+t}]$  and  $\theta$  be the specialization map that sends  $x_{s+1}, \dots, x_{s+t}$  to 0. The following statements are equivalent:

- The variables  $(x_1, \dots, x_s)$  are in Noether position with respect to the ideal  $\langle \mathbf{f} \rangle$ .
- The sequence  $(f_1, \dots, f_s, x_{s+1}, \dots, x_{s+t})$  forms a regular sequence.
- The variables  $(x_1, \dots, x_s)$  are in Noether position with respect to  $\langle \theta(f_1), \dots, \theta(f_s) \rangle$ .
- The sequence  $(\theta(f_1), \dots, \theta(f_s))$  is a regular sequence.

From a geometric point of view, when  $\mathbb{F} = \mathbb{C}$ , Noether position is strongly related to the notion of proper map below (see [7]).

**Definition 2.9.5** (Properness). Let  $V \subset \mathbb{C}^{s+t}$  be an algebraic set and  $\varphi : V \rightarrow \mathbb{C}^t$  be a polynomial morphism. The map  $\varphi$  is proper at a point  $\eta \in \mathbb{C}^t$  if there exists a neighborhood  $O$  (in the Euclidean topology) of  $\eta$  such that  $\varphi^{-1}(\overline{O})$  is bounded, where  $\overline{O}$  denotes the closure of  $O$  for the Euclidean topology over  $\mathbb{C}^t$ .

If  $\varphi$  is proper everywhere on its image, we say that the map  $\varphi$  is proper.

**Proposition 2.9.6** ([115, Proposition 3.2]). Let  $\mathbf{f} = (f_1, \dots, f_s) \in \mathbb{C}[x_1, \dots, x_{s+t}]$ . Assume that the variables  $(x_1, \dots, x_s)$  is in Noether position with respect to  $\langle \mathbf{f} \rangle$ . Then, the projection  $\pi : V(\mathbf{f}) \rightarrow \mathbb{C}^t$ ,

$$(x_1, \dots, x_{s+t}) \mapsto (x_{s+1}, \dots, x_{s+t})$$

is proper.

**Example 2.9.7.** We consider the ideal  $\langle x_1^2 + x_2^2 - 1 \rangle$ . As the polynomial  $x_1^2 + x_2^2 - 1$  is monic in  $x_1$ , the variable  $x_1$  is in Noether position with respect to this ideal.

On the other hand, the variable  $x_1$  is not in Noether position with respect to the ideal  $\langle x_1 x_2 - 1 \rangle$  as the polynomial  $x_1 x_2 - 1$  is not monic in  $x_1$ . Geometrically, this is illustrated by the fact the fiber

of the projection of  $V(x_1x_2 - 1)$  to the  $x_2$ -space tends to infinity when the value of  $x_2$  is approaching 0.

For the same reason, the variable  $x_1$  is not in Noether position with respect to  $\langle x_2x_1^2 + 2x_1 - 1 \rangle$ . In this example, the fiber over 0 of the projection of  $V(x_2x_1^2 + 2x_1 - 1)$  to the  $x_2$ -space contains a point  $(1/2, 0)$  and a point at infinity. So, this projection is not proper.

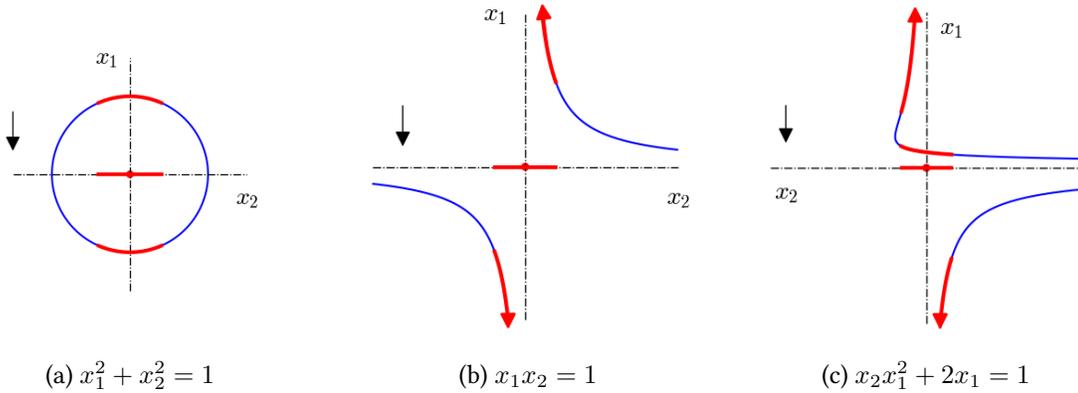


Figure 2.1: Examples for properness

# Chapter 3

## Gröbner bases

The first sections of Chapter 3 give definitions and properties of Gröbner bases, a foundation tool of symbolic computation for polynomial systems and algebraic geometry. These bases find many applications, for instance, solving the ideal membership problem, eliminating variables, computing projective closure or computing in polynomial quotient rings. The algorithmic properties of Gröbner bases will be used throughout this thesis.

We refer to the book by Cox, Little and O’Shea [48] for an introductory study on the theory of Gröbner bases and [9, Chap. 4] for the properties of Gröbner bases in zero-dimensional ideals.

### 3.1 Preliminaries on Gröbner bases

Let  $\mathbb{F}$  be a field and  $\mathbb{K}$  be an algebraic closure of  $\mathbb{F}$ . We denote by  $\mathbb{F}[\mathbf{x}]$  the polynomial ring with variables  $\mathbf{x} = (x_1, \dots, x_n)$ .

In a univariate polynomial ring  $\mathbb{F}[u]$ , Euclidean division and Euclidean algorithm are two basic operations for manipulating ideals. For instance, computing a generating polynomial of an ideal  $\langle q_1, \dots, q_s \rangle \subset \mathbb{F}[u]$  can be done by computing the greatest common divisors of  $q_1, \dots, q_s$ , an classic application of Euclidean algorithm. Testing whether a polynomial  $p \in \mathbb{F}[u]$  belongs to an ideal  $\langle q \rangle$  is done by simply checking whether  $q$  divides  $p$ .

To divide univariate polynomials, one processes successively through all the monomials in a decreasing order of degrees. However, extending the division to multivariate polynomials of  $\mathbb{F}[x_1, \dots, x_n]$  requires many more ingredients. First, one needs to define orderings among the monomials of  $\mathbb{F}[x_1, \dots, x_n]$ , which no longer depend only on the degrees as for univariate polynomials.

**Definition 3.1.1** (Monomial ordering [48, Sec. 2.2]). *An admissible monomial ordering  $\succ$  over  $\mathbb{F}[x_1, \dots, x_n]$  is a total ordering over the monomials of  $\mathbb{F}[x_1, \dots, x_n]$  satisfying the following properties:*

- $x \succ 1$  for any non-constant monomial  $x$ ;
- For any monomials  $x, y, z$  such that  $x \succ y$ , then  $zx \succ zy$ .
- Every non-empty subset of monomials of  $\mathbb{F}[x_1, \dots, x_n]$  has a smallest element.

Given a monomial ordering  $\succ$ , for a polynomial  $p \in \mathbb{F}[x_1, \dots, x_n]$ , the monomial terms of  $p$  can be ordered by  $\succ$ . Let  $c \cdot x_1^{\alpha_1} \dots x_n^{\alpha_n}$  ( $c \in \mathbb{F}$ ) be the largest monomial term of  $p$  with respect to  $\succ$ . We have the following definitions:

- The constant  $c \in \mathbb{F}$  is the leading coefficient of  $p$  with respect to  $\succ$  and is denoted by  $\text{lc}_\succ(p)$ .
- The monomial  $x_1^{\alpha_1} \dots x_n^{\alpha_n}$  is the leading monomial of  $p$  with respect to  $\succ$  and is denoted by  $\text{lm}_\succ(p)$ .
- The monomial term  $c \cdot x_1^{\alpha_1} \dots x_n^{\alpha_n}$  is the leading term of  $p$  with respect to  $\succ$  and is denoted by  $\text{lt}_\succ(p)$ .

**Example 3.1.2.** Let  $\alpha = x_1^{\alpha_1} \dots x_n^{\alpha_n}$  and  $\beta = x_1^{\beta_1} \dots x_n^{\beta_n}$  be two monomials in  $\mathbb{F}[x_1, \dots, x_n]$ . We will use the following monomial orderings:

- Lexicographic ordering  $\text{lex}(x_1 \succ \dots \succ x_n)$ :  $\alpha \succ \beta$  if and only if the left-most non-zero coefficient of  $(\alpha_1 - \beta_1, \dots, \alpha_n - \beta_n)$  is positive.
- Reverse graded lexicographic ordering  $\text{grevlex}(x_1 \succ \dots \succ x_n)$ :  $\alpha \succ \beta$  if and only if  $\deg \alpha > \deg \beta$  or  $\deg \alpha = \deg \beta$  and the right-most non-zero coefficient of  $(\alpha_1 - \beta_1, \dots, \alpha_n - \beta_n)$  is negative.
- Elimination ordering: A monomial ordering is called eliminating a subset of variables  $\mathbf{x}'$  of  $\{x_1, \dots, x_n\}$  if any  $x_i \in \{x_1, \dots, x_n\} \setminus \mathbf{x}'$  is larger than any monomial in  $\mathbf{x}'$ .

**Definition 3.1.3** (Monomial ideals). An ideal  $I \subset \mathbb{F}[x_1, \dots, x_n]$  is a monomial ideal if it can be generated by a set (not necessarily finite) of monomials of  $\mathbb{F}[x_1, \dots, x_n]$ .

The theorem below, known as Dickson's lemma, states that a monomial ideal admits a finite generating set consisting of monomials.

**Theorem 3.1.4** ([48, Chap. 2, Sec. 4, Theorem 5]). Let  $I$  be an ideal of  $\mathbb{F}[x_1, \dots, x_n]$  generated by a set  $A$  of monomials. Then  $I$  can be generated by a finite subset of  $A$ .

We fix a monomial ordering  $\succ$  over  $\mathbb{F}[x_1, \dots, x_n]$ . Let  $p \in \mathbb{F}[x_1, \dots, x_n]$  and  $(q_1, \dots, q_s)$  be a finite subset of  $\mathbb{F}[x_1, \dots, x_n]$ . With this ordering, one can already divide  $p$  by  $(q_1, \dots, q_n)$ .

Even though this division always terminates, the remainder of this division might depend on the order of the polynomials  $q_1, \dots, q_s$ . Hence, given an ideal  $I = \langle q_1, \dots, q_s \rangle$ , one may fail to test the *ideal membership* of  $p$  for  $I$  by simply choosing a "wrong" order of  $(q_1, \dots, q_s)$ . This phenomenon can be observed in Example 3.1.5 below.

**Example 3.1.5.** Let  $q_1 = x_1x_2 - 1$ ,  $q_2 = x_2^2 - 1$  with the lexicographic ordering  $x_1 \succ x_2$ . Dividing  $p = x_1x_2^2 - x_1$  by  $(q_1, q_2)$  gives

$$x_1x_2^2 - x_1 = x_2 \cdot (x_1x_2 - 1) + 0 \cdot (x_2^2 - 1) + (-x_1 + x_2).$$

However, dividing  $p$  by  $(q_2, q_1)$ , we have

$$x_1x_2^2 - x_1 = x_1 \cdot (x_2^2 - 1) + 0 \cdot (x_1x_2 - 1) + 0.$$

From the second division, we know that  $p \in \langle q_1, q_2 \rangle$  while this can not be obtained from the first division.

To overcome this weakness, in 1976, Buchberger introduced the notion of Gröbner bases [30] which then becomes one of the foundations of computer algebra. In what follows, we recall the definition and some preliminary results of Gröbner bases.

Further, we fix a monomial ordering  $\succ$  over  $\mathbb{F}[x_1, \dots, x_n]$ . Let  $I$  be an ideal of  $\mathbb{F}[x_1, \dots, x_n]$ . The initial ideal of  $I$  with respect to the ordering  $\succ$  is the ideal

$$\langle \text{lm}_\succ(p) \mid p \in I \rangle.$$

Given a set  $(f_1, \dots, f_s)$  of generators of  $I$ , in general, the ideal

$$\langle \text{lm}_\succ(f_i) \mid 1 \leq i \leq s \rangle$$

is not equal to the initial ideal of  $I$ .

**Example 3.1.6.** Let  $f_1 = x_1^2 + 2x_1x_2 + x_2^2 + 2x_1 + 1$  and  $f_2 = x_1^2 + 2x_1x_2 + 2x_2^2 + x_1 + x_2$ . Taking the lexicographic ordering  $x_1 \succ x_2$ , the leading monomials of  $f_1$  and  $f_2$  are both  $x_1^2$ . On the other hand, since

$$f_3 = x_1 - x_2^2 - x_2 + 1 = f_1 - f_2$$

belongs to  $\langle f_1, f_2 \rangle$ , the initial ideal of  $\langle f_1, f_2 \rangle$  contains the leading monomial of  $f_3$  which is  $x_1$ .

Definition 3.1.7 below defines the Gröbner bases of an ideal  $I \subset \mathbb{F}[x_1, \dots, x_n]$ , which are finite generating sets of  $I$  whose leading monomials also generate the initial ideal of  $I$ .

**Definition 3.1.7** (Gröbner bases). Let  $I$  be an ideal of  $\mathbb{F}[x_1, \dots, x_n]$ . A Gröbner basis  $G$  of  $I$  with respect to the ordering  $\succ$  is a finite subset of  $I$  such that the set of leading monomials  $\{\text{lm}_\succ(g) \mid g \in G\}$  generates the initial ideal  $\langle \text{lm}_\succ(p) \mid p \in I \rangle$ .

**Proposition 3.1.8** ([48, Ch. 2, Sec. 5, Cor. 6]). Let  $G$  be a Gröbner basis of  $I$  with respect to the ordering  $\succ$ . Then,  $G$  is a generating set of the ideal  $I$ .

Also in [30], Buchberger presented a criterion to decide whether a set of polynomials is a Gröbner basis. From his criterion, he derived the first algorithm to compute Gröbner bases. These algorithmic results are based on the construction of  $S$ -polynomials and their properties.

**Definition 3.1.9** ( $S$ -polynomials). Let  $f, g \in \mathbb{F}[x_1, \dots, x_n]$ . The  $S$ -polynomial of  $f$  and  $g$  is defined as

$$S(f, g) = \text{lcm}(\text{lm}_\succ(f), \text{lm}_\succ(g)) \left( \frac{f}{\text{lt}_\succ(f)} - \frac{g}{\text{lt}_\succ(g)} \right).$$

**Proposition 3.1.10** (Buchberger criterion, [48, Chap. 2, Sec. 6, Theorem 6]). Let  $I$  be a polynomial ideal. Then a basis  $G = \{g_1, \dots, g_s\}$  of  $I$  is a Gröbner basis of  $I$  if and only if for all pairs  $i \neq j$ , the remainder on division of  $S(g_i, g_j)$  by  $G$  listed in some order is zero.

Through an iterative procedure, Buchberger algorithm tries to discover new initial terms by adding more polynomials to the generating set. For each pair  $(f, g)$  of polynomials in the current basis, the algorithm computes the  $S$ -polynomial  $S(f, g)$  and reduces it by the current basis. If the remainder of this division is not zero, this remainder is added as a new element to the basis.

The algorithm terminates when no new polynomial can be added and the output is a Gröbner basis of the input ideal with respect to the considered ordering (see [30] or [48, Chap. 2, Sec. 7, Theorem 2]). Note that the basic version of Buchberger algorithm we mention above is mostly of theoretical interest. For practical computations, many additional criteria and optimizations have been proposed to improve its performance. We will discuss about these variants of Buchberger algorithm and more efficient algorithms for computing Gröbner bases in Subsection 3.4.

An important property of Gröbner bases is the uniqueness of remainders of polynomial divisions, which provides an algorithm for the ideal membership problem (Proposition 3.1.13).

**Proposition 3.1.11** ([48, Chap. 2, Sec. 6, Prop. 1]). *Let  $I$  be an ideal of  $\mathbb{F}[x_1, \dots, x_n]$  and  $G$  be a Gröbner basis of  $I$  with respect to some ordering  $\succ$ . Given  $p \in \mathbb{F}[x_1, \dots, x_n]$ , the remainder of the division of  $p$  by  $G$  using the monomial ordering  $\succ$  is uniquely defined. It is called the normal form of  $p$  with respect to  $G$  and is denoted by  $\text{NF}_G(p)$ .*

We continue with Example 3.1.5.

**Example 3.1.12.** *Let  $g_1 = x_1 - x_2$ ,  $g_2 = x_2^2 - 1$  be two polynomials in  $\mathbb{C}[x_1, x_2]$ . Using Buchberger criterion, one can verify that  $\{g_1, g_2\}$  is a Gröbner basis of the ideal  $\langle x_1x_2 - 1, x_2^2 - 1 \rangle$  with respect to the lexicographic ordering  $x_1 \succ x_2$ .*

*Dividing  $p = x_1x_2^2 - x_1$  respectively by  $(g_1, g_2)$  and  $(g_2, g_1)$  we obtain*

$$\begin{aligned} x_1x_2^2 - x_1 &= (x_2^2 - 1) \cdot (x_1 - x_2) + x_2 \cdot (x_2^2 - 1) + 0, \\ &= x_1 \cdot (x_2^2 - 1) + 0 \cdot (x_1 - x_2) + 0. \end{aligned}$$

*The remainders are 0 in the both cases.*

**Proposition 3.1.13** ([48, Chap. 2, Sec. 6, Cor. 2]). *Let  $I$  be an ideal of  $\mathbb{F}[x_1, \dots, x_n]$  and  $G$  be a Gröbner basis of  $I$  with respect to any monomial ordering. Then  $p \in I$  if and only if  $\text{NF}_G(p) = 0$ .*

**Definition 3.1.14.** *Let  $I$  be an ideal of  $\mathbb{F}[x_1, \dots, x_n]$  and  $G$  be a Gröbner basis  $G$  of  $I$  with respect to some ordering  $\succ$ . We have the following definitions:*

- $G$  is called *minimal* if for every pair  $g_i, g_j \in G$ ,  $\text{lt}_\succ(g_i)$  does not divide  $\text{lt}_\succ(g_j)$ .
- $G$  is called *reduced* if for any polynomial  $g \in G$ ,  $g$  is monic and no monomial of  $g$  lies in  $\langle \text{lt}_\succ(g') \mid g' \in G \setminus \{g\} \rangle$ .

**Proposition 3.1.15** ([48, Chap. 2, Sec. 5, Cor. 6]). *Let  $I$  be an ideal of  $\mathbb{F}[x_1, \dots, x_n]$ . For any monomial ordering  $\succ$ , there exists a unique reduced Gröbner basis  $G$  of  $I$  with respect to  $\succ$ .*

Proposition 3.1.15 gives an algorithm for testing equality between two ideals. An immediate corollary is a test whether  $V_{\mathbb{K}}(I)$  is empty.

**Proposition 3.1.16.** *Let  $I$  be an ideal of  $\mathbb{F}[x_1, \dots, x_n]$ . Then,  $V_{\mathbb{K}}(I) = \emptyset$ , i.e.,  $I = \langle 1 \rangle$ , if and only if any Gröbner basis of  $I$  contains a non-zero element of  $\mathbb{F}$ .*

## 3.2 Algebraic elimination using Gröbner bases

We already see from Propositions 3.1.13 and 3.1.15 that Gröbner bases allow one to test the membership of a polynomial with respect to an ideal or the equality of two polynomial ideals. In this section, we illustrate how to use Gröbner bases to carry out algebraic elimination.

**Theorem 3.2.1** ([48, Chap. 3, Sec. 2, Theorem 2]). *Let  $V \subset \mathbb{K}^n$  be an algebraic set and  $\pi$  be the projection*

$$(x_1, \dots, x_n) \mapsto (x_{k+1}, \dots, x_n).$$

*Then, the algebraic set defined by*

$$I \cap \mathbb{K}[x_{k+1}, \dots, x_n]$$

*is the Zariski closure of  $\pi(V)$ .*

**Theorem 3.2.2** (Elimination theorem, [48, Chap. 3, Sec. 1, Theorem 2]). *Let  $I$  be an ideal of  $\mathbb{F}[x_1, \dots, x_n]$  and  $G$  be a Gröbner basis of  $I$  with respect to an ordering eliminating the variables  $x_1, \dots, x_k$ .*

*Then, we have*

$$I \cap \mathbb{F}[x_{k+1}, \dots, x_n] = \langle G \cap \mathbb{F}[x_{k+1}, \dots, x_n] \rangle.$$

Theorems 3.2.1 and 3.2.2 provide an algorithm for computing the Zariski closures of projections of algebraic sets. One computes a Gröbner basis  $G$  of  $I(V)$  with respect to an ordering eliminating the variables  $x_1, \dots, x_k$  and takes all the elements of  $G$  that contains only the variables  $x_{k+1}, \dots, x_n$ .

**Example 3.2.3.** *We compute the critical values of the restriction of  $\pi : (x_1, x_2) \mapsto x_1$  to the sphere  $\mathcal{V} \subset \mathbb{C}^2$  defined by*

$$x_1^2 + x_2^2 = 1.$$

*By Jacobian criterion, the set of critical points  $\text{crit}(\pi, \mathcal{V})$  can be defined by*

$$x_1^2 + x_2^2 = x_2 = 0.$$

*The reduced Gröbner basis of  $\langle x_1^2 + x_2^2, x_2 \rangle$  with respect to the ordering  $\text{lex}(x_2 \succ x_1)$  is*

$$\{x_2, x_1^2 - 1\}.$$

Thus, we obtain the critical values by taking the intersection

$$\langle x_2, x_1^2 - 1 \rangle \cap \mathbb{C}[x_1] = \langle x_1^2 - 1 \rangle,$$

which gives  $x_1 = 1$  and  $x_1 = -1$ .

Algebraic elimination can also be used to compute saturation ideals. Recall that the saturation ideal  $I : J^\infty$  where  $I, J$  are two ideals of  $\mathbb{F}[x_1, \dots, x_n]$  is defined as

$$I : J^\infty = \{f \in \mathbb{F}[x_1, \dots, x_n] \mid \exists k \in \mathbb{Z}_+ \text{ such that } fJ^k \subset I\}.$$

Geometrically, the algebraic set defined by  $I : J^\infty$  is the Zariski closure of  $V(I) \setminus V(J)$ . Proposition 3.2.4 below deduces an algorithm for computing the saturation ideal when  $J$  is generated by one polynomial.

**Proposition 3.2.4** ([48, Chap. 4 Sec. 4 Theorem 14.]). *Let  $I = \langle f_1, \dots, f_s \rangle$  be an ideal of  $\mathbb{F}[x_1, \dots, x_n]$  and  $g \in \mathbb{F}[x_1, \dots, x_n]$ . Then,*

$$I : g^\infty = \langle f_1, \dots, f_s, \ell \cdot g - 1 \rangle \cap \mathbb{F}[x_1, \dots, x_n].$$

There is also an algorithm for computing saturation ideals due to Bayer (see [12] or [56, Chap. 18]) which is known to be faster than the algorithm above.

Another important application of Gröbner bases comes from the following proposition.

**Proposition 3.2.5** ([48, Chap. 5, Sec. 3, Prop. 4]). *Let  $I \subset \mathbb{F}[x_1, \dots, x_n]$  be an ideal and  $G$  be a Gröbner basis of  $I$  for some monomial ordering. The set of monomials in  $\mathbb{F}[x_1, \dots, x_n]$  which are not reducible by  $G$  forms a basis of the  $\mathbb{F}$ -vector space  $\mathbb{F}[x_1, \dots, x_n]/I$ .*

From this property, one derives an algorithm for computing representatives of elements in the quotient ring  $\mathbb{F}[x_1, \dots, x_n]/I$ . These representatives allow explicitly arithmetic computations over  $\mathbb{F}[x_1, \dots, x_n]/I$ , especially for a zero-dimensional ideal  $I$  that we will see in the next section.

### 3.3 Gröbner bases and zero-dimensional ideals

Let  $\mathbb{F}$  be a field and  $\mathbb{K}$  be an algebraic closure of  $\mathbb{F}$ .

An ideal  $I$  of  $\mathbb{F}[x_1, \dots, x_n]$  is said to be zero-dimensional if the algebraic set  $V_{\mathbb{K}}(I) \subset \mathbb{K}^n$  is finite and non-empty. The following proposition characterizes the zero-dimensional ideals.

**Theorem 3.3.1** ([48, Sec. 5.3, Theorem 6]). *Let  $I$  be an ideal of  $\mathbb{F}[x_1, \dots, x_n]$ . The following statements are equivalent:*

- *The ideal  $I$  is zero-dimensional.*
- *The quotient ring  $\mathbb{F}[x_1, \dots, x_n]/I$  is a  $\mathbb{F}$ -vector space of finite positive dimension.*

- For every  $x_i$ , there exists a univariate polynomial  $p_i \in \mathbb{F}[x_i]$  such that  $p_i \in I$ .

The dimension of  $\mathbb{F}[x_1, \dots, x_n]/I$  as  $\mathbb{F}$ -vector space is called the algebraic degree of  $I$ , denoted by  $\deg(I)$ .

Further in this section, we let  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{K}[x_1, \dots, x_n]$  be a polynomial sequence generating a zero-dimensional ideal  $I = \langle \mathbf{f} \rangle$  (hence,  $s \geq n$ ). We denote by  $A$  the quotient ring  $\mathbb{K}[x_1, \dots, x_n]/I$ .

The solution set  $V_{\mathbb{K}}(I)$  is a finite set. For each  $\eta = (\eta_1, \dots, \eta_n) \in V_{\mathbb{K}}(I)$ , recall that the ideal of elements of  $\mathbb{K}[x_1, \dots, x_n]$  vanishing at  $\eta$  is the maximal ideal

$$\langle x_1 - \eta_1, \dots, x_n - \eta_n \rangle \subset \mathbb{K}[x_1, \dots, x_n].$$

We denote this ideal by  $I_\eta$ .

The cardinality of  $V_{\mathbb{K}}(I)$  is at most  $\deg(I)$ , and it is exactly  $\deg(I)$  if points are counted with the following notion of multiplicity.

We define now a notion of multiplicity for each point of  $V_{\mathbb{K}}(I)$ .

**Definition 3.3.2.** Let  $\eta \in V_{\mathbb{K}}(I)$ . The image  $\bar{I}_\eta$  of  $I_\eta$  in  $A$  is also a maximal ideal of  $A$ .

The localization of  $A$  at  $\bar{I}_\eta$  is a local ring, denoted by  $A_\eta$ .

The structure of  $A$  can be read off from the local rings  $A_\eta$ .

**Theorem 3.3.3** ([9, Theorem 4.95]). We have the following ring isomorphism:

$$A \simeq \prod_{\eta \in V_{\mathbb{K}}(I)} A_\eta.$$

As a consequence, the local ring  $A_\eta$  is a  $\mathbb{K}$ -vector space of finite dimension.

**Definition 3.3.4.** We define the multiplicity of  $\eta$ , denoted by  $\mu(\eta)$ , as the dimension of  $A_\eta$  as  $\mathbb{K}$ -vector space.

A point  $\eta \in V_{\mathbb{K}}(I)$  is a singular zero of  $\mathbf{f}$  if the rank of the Jacobian matrix

$$\begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_s}{\partial x_1} & \dots & \frac{\partial f_s}{\partial x_n} \end{bmatrix}$$

at  $\eta$  is at most  $n - 1$ . Otherwise, if this matrix has rank  $n$  at  $\eta$ ,  $\eta$  is a non-singular zero.

The following proposition states that the singular solutions of  $I$  are the ones with multiplicities greater than 1.

**Proposition 3.3.5** ([9, Prop. 4.96]). Let  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{K}[x_1, \dots, x_n]$  be a zero-dimensional system and  $\eta = (\eta_1, \dots, \eta_n)$  a zero of  $\mathbf{f}$  in  $\mathbb{K}^n$ . The ideal  $\langle x_1 - \eta_1, \dots, x_n - \eta_n \rangle$  is denoted by  $I_\eta$ . Then, the following are equivalent:

- $\eta$  is a non-singular zero of  $\mathbf{f}$ .
- The multiplicity of  $\eta$  is 1.
- $I_\eta \subset \langle \mathbf{f} \rangle + I_\eta^2$ .

The theorem below, known as Bézout theorem, gives a bound on the number of complex solutions (counted with multiplicities) of a zero-dimensional ideal through the degrees of its defining equations.

**Theorem 3.3.6** (Bézout theorem, [101, Theorem 1]). *Let  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{K}[x_1, \dots, x_n]$  be a zero-dimensional system. The dimension of the quotient ring  $\mathbb{K}[x_1, \dots, x_n]/I$  as  $\mathbb{K}$ -vector space is bounded by  $\deg(f_1) \dots \deg(f_s)$ .*

The rest of this section is used to explain how to compute over zero-dimensional ideals. We begin with the computation over the quotient ring  $\mathbb{F}[x_1, \dots, x_n]/I$  with the knowledge of a Gröbner basis of  $I$ .

Fixing a monomial ordering  $\succ$  over  $\mathbb{F}[x_1, \dots, x_n]$ , let  $G$  be the reduced Gröbner basis of  $I$  with respect to  $\succ$  and  $B$  be the set of monomials in  $\mathbb{F}[x_1, \dots, x_n]$  which are not reducible by  $G$ . By Proposition 3.2.5,  $B$  is a basis of  $\mathbb{F}[x_1, \dots, x_n]/I$  as a  $\mathbb{F}$ -vector space.

For any  $p \in \mathbb{F}[x_1, \dots, x_n]$ , the normal form of  $p$  by  $G$  is a linear combination of elements of  $B$  with coefficients in  $\mathbb{F}$ . This normal form can be interpreted as the image of  $p$  in  $\mathbb{F}[x_1, \dots, x_n]/I$ . Therefore, the operations in the quotient ring  $\mathbb{F}[x_1, \dots, x_n]/I$  such as vector additions or scalar multiplications can be computed explicitly using the normal form reduction.

The finite-dimensional vector space structure of  $\mathbb{F}[x_1, \dots, x_n]$  and the computation powered by Gröbner bases are the main ingredients for the ordering change algorithm FGLM [67] (see Subsection 3.4.4) and the construction of Hermite matrix (see Subsection 4.4.3).

Now we discuss how to obtain explicitly the solutions of a given zero-dimensional ideal  $I$ . Usually, a Gröbner basis with respect to lexicographic ordering provides a triangular description of  $I$ , which allows one to retrieve the solutions of  $I$  coordinate by coordinate. More specifically, we define the notion *shape position*.

**Definition 3.3.7** (Shape position). *A zero-dimensional ideal  $I \subset \mathbb{F}[x_1, \dots, x_n]$  is in shape position if the reduced Gröbner basis of  $I$  with respect to the lexicographic ordering  $x_1 \succ \dots \succ x_n$  has the following form*

$$\{x_1 - g_1, \dots, x_{n-1} - g_{n-1}, \dots, g_n\},$$

where  $g_1, \dots, g_n$  lie in  $\mathbb{F}[x_n]$ .

The following proposition, known as *Shape lemma*, claims that substituting  $x_n$  by a generic linear form will help bring any ideal in shape position.

**Proposition 3.3.8** ([14]). *Let  $I$  be a radical zero-dimensional ideal of  $\mathbb{F}[x_1, \dots, x_n]$ . There exists a non-empty Zariski open subset  $\mathcal{A}$  of  $\mathbb{K}^n$  such that for any  $(\lambda_1, \dots, \lambda_n) \in \mathcal{A}$ , the ideal*

$$\{p(x_1, \dots, x_{n-1}, \lambda_1 x_1 + \dots \lambda_n x_n) \mid p \in I\}$$

is in shape position with respect to  $x_1, \dots, x_n$ .

The data structure we use to represent finite algebraic sets is the following *zero-dimensional parametrization*, which uses also the idea of projecting on a generic linear form.

**Definition 3.3.9.** A zero-dimensional parametrization  $\mathcal{R}$  of coefficients in  $\mathbb{Q}$  over  $\mathbb{C}^n$  consists of

- A square-free polynomial  $w \in \mathbb{Q}[u]$  where  $u$  is a new variable;
- A sequence  $(\lambda_1, \dots, \lambda_n) \in \mathbb{Q}^n$  such that

$$u \cdot w' = \sum_{i=1}^n \lambda_i \cdot v_i \pmod{w};$$

- A sequence of polynomials  $(v_1, \dots, v_n)$  in  $\mathbb{Q}[u]$  with  $\deg(v_i) < \deg(w)$ .

The solution set of  $\mathcal{R}$ , denoted by  $Z(\mathcal{R})$ , is the following finite set

$$Z(\mathcal{R}) = \left\{ \left( \frac{v_1(u)}{w'(u)}, \dots, \frac{v_n(u)}{w'(u)} \right) \mid w(u) = 0 \right\}.$$

A zero-dimensional algebraic set  $V \subset \mathbb{C}^n$  is represented by a zero-dimensional parametrization  $\mathcal{R}$  if and only if  $V$  coincides with  $Z(\mathcal{R})$ .

Given a polynomial sequence  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{Q}[x_1, \dots, x_n]$  such that  $V_{\mathbb{C}}(\mathbf{f})$  is finite, the knowledge of a zero-dimensional parametrization of  $V_{\mathbb{C}}(\mathbf{f})$  allows one to extract numerical values of the solutions of  $\mathbf{f}$  up to arbitrary precision using root isolating algorithms ([168, 122]).

To compute such parametrizations, we refer to algorithms such as the rational univariate representation [166] or the geometric resolution [87]. While the first algorithm relies on computing a Gröbner basis of the ideal  $\langle \mathbf{f} \rangle$ , the second one, which is probabilistic, uses a process of incrementally lifting and intersecting curves.

Note that it is possible to retrieve a polynomial parametrization by inverting the derivative  $w'$  modulo  $w$ . Still, the parametrization with  $w'$  as denominator is known to be better for practical computations as it usually involves coefficients with smaller bit size (see [49]).

**Example 3.3.10.** We consider the system

$$x_1^2 + 2x_2^2 + 2x_1 + 1 = 2x_1x_2 + 2x_2^2 + x_1 + x_2 = 0.$$

A zero-dimensional parametrization with  $(\lambda_1, \lambda_2) = (1, 2)$  is

$$(w, v_1, v_2) = (6u^4 + 20u^3 + 13u^2 - 10u + 9, -2(8u^3 + 17u^2 - 7u - 4), -2(u^3 - 2u^2 - 4u + 11)).$$

### 3.4 On the computation of Gröbner bases

Let  $\mathbf{f} = (f_1, \dots, f_s)$  be a polynomial sequence in  $\mathbb{F}[x_1, \dots, x_n]$  with a monomial ordering  $\succ$ . In this section, we discuss about algorithms that compute Gröbner bases of  $\langle \mathbf{f} \rangle$  with respect to  $\succ$ .

### 3.4.1 Buchberger algorithm's drawbacks

As we mentioned in Section 3.1, the first algorithm for computing Gröbner bases is due to Buchberger [30]. This algorithm is based on two operations:

- Choosing a pair  $(f, g)$  from the current basis and constructing the  $S$ -polynomial  $S(f, g)$ ;
- Reducing  $S(f, g)$  to the current basis.

It is quickly observed that, most of the reductions of  $S$ -polynomials result in zeros and do not play any further role in the computation. This leads to inefficiency of this algorithm in practice.

On the other hand, from the classical Buchberger algorithm, the pair  $(f, g)$  is chosen freely among the current generating set. This choice affects the run of algorithms.

Therefore, the common ideas for improving Buchberger algorithm are to mitigate the two weaknesses above, i.e., to avoid as much as possible the useless reductions to zero and to have a good strategy of choosing critical pairs. This first problem was addressed in [31], where Buchberger introduced two extra criteria for his algorithm to filter out some pairs leading to reductions to zero. Several strategies for the choice of critical pairs are also proposed, for instance, the *sugar strategy* [79]. The variants of Buchberger algorithm with these improvements are implemented in Macaulay2, Singular, Magma or Maple.

F4/F5 algorithms that we discuss in the next subsection also address these two drawbacks.

### 3.4.2 F4/F5 algorithms

In 1999, Faugère presented F4 algorithm [59], which uses the same mathematical principles as Buchberger algorithm but with a linear algebra approach. This algorithm constructs a matrix indexed by monomials up to some degree and carry out many  $S$ -polynomials reductions at once by an echelon form reduction of the mentioned matrix. This strategy actually avoids the choice of critical pairs by processing many of them at once, which surprisingly is very efficient. Showing an impressive practical behavior, F4 algorithm is widely implemented for applications. Its implementations are available in libraries such as FGB [66], MSOLVE [17] or compute algebra systems like Magma or Maple.

Later on in [60], Faugère introduced the signatures of polynomials for avoiding reductions to zero. This notion leads to a new efficient algorithm, named F5, for computing Gröbner bases. Especially, when the input system forms a regular sequence (which is true for generic systems), F5 algorithm does not perform any reduction to zero. In practice, it successfully solved a set of 80 dense polynomials in 80 variables over some finite field [62].

### 3.4.3 Complexity issues

We now discuss the complexity of computing Gröbner bases.

In a series of works in [146, 80, 151], the worst case complexity of Gröbner bases is proved to lie in  $2^{2^{O(n)}}$  where  $n$  is the number of variables. However, this complexity is only obtained for extremely rare systems which are constructed on purpose and is not reached by generic systems.

In practice, it has been observed that the actual behavior of Gröbner basis implementations can be quite efficient. This motivates the studies of complexity of Gröbner bases, which lead to nice results for useful special classes of polynomial systems.

It is worth noting that different orderings behave differently in computation. While lexicographic orderings provide an explicit triangular description of a polynomial system, the graded reversed lexicographic (grevlex) ordering generally yields Gröbner bases of smaller degrees and coefficients. Algorithms such as F4/F5 are much more efficient when performing with those orderings. Therefore, the complexity analysis is mostly carried out for grevlex orderings.

Besides, the complexity of the previous algorithms is difficult to estimate. It mainly depends on the number of critical pairs remained to be processed at any given step. However, the complexity of these algorithms can be bounded by studying their matrix variants, replacing  $S$ -polynomials and critical pairs with the construction and reduction of Macaulay matrices.

Matrix-F5 algorithm is designed in [7] as a variant of F5 algorithm which is well-suited for complexity analysis. This algorithm takes an additional parameter along with the system, a degree  $D_{\max}$  at which to stop. Formally, what the algorithm computes is a Gröbner basis truncated to  $D_{\max}$ , which is a set of polynomials containing the polynomials of the reduced Gröbner basis of  $I$  of degree at most  $D_{\max}$ .

**Proposition 3.4.1** ([7, Prop. 1]). *Let  $(f_1, \dots, f_s)$  be a system of homogeneous polynomials in  $\mathbb{F}[x_1, \dots, x_n]$ . The number of operations in  $\mathbb{F}$  required to compute a Gröbner basis of  $\langle f_1, \dots, f_s \rangle$  for grevlex ordering up to degree  $D_{\max}$  is bounded by*

$$O\left(sD_{\max} \binom{n + D_{\max} - 1}{D_{\max}}^\omega\right)$$

where  $\omega$  is the exponent of matrix multiplication.

This complexity statement immediately leads to the question of finding the required value of  $D_{\max}$  to obtain a complete Gröbner basis. The highest degree appearing in a Gröbner basis  $G$  is called the *degree of regularity* of  $G$ .

For zero-dimensional ideals, Lazard has shown in [133] that after a generic linear change of coordinates, the degree of regularity  $D_{\text{reg}}$  for computing a Gröbner basis of grevlex ordering is bounded by

$$D_{\text{reg}} \leq \sum_{i=1}^s (\deg(f_i) - 1) + 1.$$

This bound is known as *Macaulay's bound*.

In [7], this bound is extended to positive dimensional systems  $(f_1, \dots, f_s)$  under some genericity assumptions to obtain a *simply* exponential complexity in the number of variables  $n$  for computing Gröbner bases.

**Theorem 3.4.2** ([7, Prop. 1, Theorem 12]). *Let  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{F}[x_1, \dots, x_n]$  be a system of homogeneous polynomials whose degrees are uniformly bounded by  $D$ .*

Assume that  $\mathbf{f}$  is in Noether position with respect to the variables  $x_{s+1}, \dots, x_n$ . Then, Matrix-F5 algorithm computes a Gröbner basis of  $\langle \mathbf{f} \rangle$  for grevlex ordering within

$$O\left(s^2 D \binom{n + s(D-1)}{s(D-1) + 1}^\omega\right)$$

arithmetic operations in  $\mathbb{F}$ .

Note also that, more complicated analyses are also given in [7] to refine the exponent  $\omega$  in the above complexity. However, these details are not relevant to our thesis so we satisfy with the complexity established in Theorem 3.4.2.

### 3.4.4 Change of monomial ordering algorithms

In this subsection, we present a strategy for computing Gröbner bases through changing the monomial orderings.

As explained in the previous subsection, the grevlex orderings are more suitable for computations while a lex Gröbner basis can provide more geometric information. A natural idea is to compute a Gröbner basis with respect to the grevlex ordering and then, from this basis, one calls to an ordering change algorithm to obtain the lex basis. For *zero-dimensional ideals*, this can be efficiently done by FGLM algorithm, named after Faugère, Gianni, Lazard and Mora [67], by exploiting the finite-dimensional vector-space structure of the quotient ring.

Let  $(f_1, \dots, f_s) \subset \mathbb{F}[x_1, \dots, x_n]$  that generates a zero-dimensional ideal  $I$  and  $\succ_1, \succ_2$  be two monomial orderings in  $\mathbb{F}[x_1, \dots, x_n]$ . From a Gröbner basis  $G_1$  of  $I$  with respect to  $\succ_1$ , FGLM algorithm computes a Gröbner basis  $G_2$  with respect to  $\succ_2$ .

More specifically, it proceeds as follows.

- Use  $G_1$  to compute the basis  $B_1$  of  $\mathbb{F}[x_1, \dots, x_n]/I$  with respect to  $\succ_1$ ;
- Compute the multiplication matrices of  $x_1, \dots, x_n$  with respect to  $B_1$ ;
- Compute the staircase of  $I$  with respect to  $\succ_2$  and derive  $G_2$ .

The two first steps are performed following the explanation in Section 3.3. The last step is carried out by constructing a basis  $B_2$  of  $\mathbb{F}[x_1, \dots, x_n]/I$  using linear algebra in  $\mathbb{F}[x_1, \dots, x_n]/I$ .

Let  $D$  be a bound on the degree of  $f_i$  and  $\delta$  be the dimension of the quotient ring  $\mathbb{F}[x_1, \dots, x_n]$  as a  $\mathbb{F}$ -vector space. The arithmetic complexity of the original FGLM algorithm is bounded by

$$O(n\delta^3).$$

A faster variant of FGLM algorithm is proposed by Faugère and Mou [68]. This new algorithm takes advantage of the sparsity of multiplication maps  $\mathcal{L}_{x_i}$  to obtain a complexity

$$O\left(\sqrt{\frac{6}{n\pi}} \delta^{2+\frac{n-1}{n}}\right).$$

Recently in [58], this variant is improved using fast linear algebra to run within

$$O((Dn^{\omega+1} + \log_2(\delta))\delta^\omega)$$

arithmetic operations in  $\mathbb{F}$  for a regular sequence  $(f_1, \dots, f_n) \subset \mathbb{F}[x_1, \dots, x_n]$ . Combining with the computation of a grevlex Gröbner basis, [58] results in a faster algorithm for solving any zero-dimensional system  $\mathbf{f} = (f_1, \dots, f_n)$  which generates a radical ideal and forms an affine regular sequence in  $\mathbb{Q}[x_1, \dots, x_n]$ . This algorithm has an arithmetic complexity lying within

$$O(nD^{\omega n} + (Dn^{\omega+1} + \log_2(\delta))\delta^\omega).$$

# Chapter 4

## Basic notions of real algebraic geometry

Real algebraic geometry studies the solutions of polynomial systems of equations and inequalities over a real field. It has many different properties compared to the classic algebraic geometry that focuses on the algebraically closed field.

Chapter 4 presents some basic definitions and results in the real algebraic geometry which are used in the thesis. Most of the results are presented in the book by Basu, Pollack and Roy [9].

### 4.1 Real fields

A real field is a field over which one can define an ordering of elements. This type of fields is a generalization of the field of real numbers  $\mathbb{R}$  and shares many similar properties with  $\mathbb{R}$ .

We first give the precise definition of a real field.

**Definition 4.1.1** (Real field). *An ordering of a field  $R$  is a total order relation  $\leq$  satisfying:*

- $x \leq y$  then  $x + z \leq y + z$ ,
- $0 \leq x, y$  then  $0 \leq xy$ .

*A field  $R$  is a real field if it can be equipped with an ordering.*

**Example 4.1.2.** *The fields  $\mathbb{Q}$  and  $\mathbb{R}$  are real fields with usual order over  $\mathbb{R}$ .*

Below we give some characterization of real fields.

**Proposition 4.1.3** ([19, Theorem 1.18]). *The following statements are equivalent*

- $R$  is a real field.
- $-1$  is not a sum of squares of elements of  $R$ .
- For any  $f_1, \dots, f_s$  such that

$$f_1^2 + \dots + f_s^2 = 0,$$

*then  $f_1 = \dots = f_s = 0$ .*

**Definition 4.1.4** (Real closed field [9, Theorem 2.17]). *A real field  $R$  is called a real closed field if the extension  $R[u]/\langle u^2 + 1 \rangle$  is an algebraically closed field.*

**Example 4.1.5.** *We name a few examples for real closed fields:*

- The field  $\mathbb{R}$  is a real closed field.
- A real number is called algebraic if it is a root of some univariate polynomial with integer coefficients. The field of real algebraic numbers, denoted by  $\mathbb{R}_{\text{alg}}$ , is a real closed field.

One can define the intermediate value property over an arbitrary real field which is similar to the classical one in  $\mathbb{R}$ . It provides a useful tool for proving properties of continuous maps.

**Definition 4.1.6** (Intermediate value property). *A field  $R$  has the intermediate value property if  $R$  is a real field such that, for any  $p \in R[u]$ , if there exist  $a, b \in R$  such that  $p(a) \cdot p(b) < 0$ , then there exists  $c \in (a, b)$  such that  $p(c) = 0$ .*

**Theorem 4.1.7** ([9, Theorem 2.17]). *A real field  $R$  is closed if and only if it has the intermediate value property.*

## 4.2 Semi-algebraic sets

Let  $R$  be a real closed field. We study the solution sets of polynomial systems consisting of equations and inequalities with coefficients in  $R$ ; these sets are called semi-algebraic sets.

**Definition 4.2.1** (Semi-algebraic sets). *A semi-algebraic formula  $\Phi$  defined over  $R[x_1, \dots, x_n]$  is a logic formula of the form*

$$\bigvee_{i=1}^m \bigwedge_{j=1}^{n_i} f_{i,j} \bullet_{i,j} 0,$$

where  $f_{i,j} \in R[x_1, \dots, x_n]$  and  $\bullet_{i,j} \in \{>, =\}$ . The polynomials  $f_{i,j}$  are called atoms of  $\Phi$ .

A subset  $S$  of  $R^n$  is a semi-algebraic set if there exists a semi-algebraic formula  $\Phi$  defined over  $R[x_1, \dots, x_n]$  such that

$$S = \{\eta \in R^n \mid \Phi(\eta) \text{ is true}\}.$$

When a semi-algebraic set can be defined by equations only, it is called a real algebraic set.

Note that every real algebraic set defined by  $f_1 = \dots = f_s = 0$  can also be defined by only one equation  $f_1^2 + \dots + f_s^2 = 0$ .

**Definition 4.2.2** (Semi-algebraic maps). *Let  $S \subset R^n$  and  $S' \subset R^p$  be two semi-algebraic sets and  $\varphi$  be a map from  $S$  to  $S'$ . The graph of  $\varphi$  is defined as the subset*

$$\{(\eta, \varphi(\eta)) \mid \eta \in S\}$$

of  $S \times S'$ . The map  $\varphi$  is called semi-algebraic if its graph is a semi-algebraic set of  $R^n \times R^p$ .

The property below is easily deduced from the definition of semi-algebraic maps.

**Lemma 4.2.3.** *Let  $\varphi : S \rightarrow S'$  be a semi-algebraic map which is bijective. Then  $\varphi^{-1} : S' \rightarrow S$  is also semi-algebraic.*

From the definition, the set of semi-algebraic sets is stable under finite union, intersection and complement. On the other hand, unlike algebraic sets whose projections are not necessarily algebraic, the image of a semi-algebraic set by a semi-algebraic map (including projections) is also a semi-algebraic set. This important result is due to Tarski [192] and Seidenberg [182].

**Theorem 4.2.4** (Tarski-Seidenberg theorem, [9, Theorem 2.98]). *Given two semi-algebraic sets  $S \subset R^n$ ,  $S' \subset R^p$  and a semi-algebraic map  $\varphi : S \rightarrow S'$ , then  $\varphi(S) \subset R^p$  is semi-algebraic.*

Given a semi-algebraic formula  $\Phi$  whose atoms are in  $R[\mathbf{x}, \mathbf{y}]$  where  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{y} = (y_1, \dots, y_t)$ , this semi-algebraic formula defines a semi-algebraic set  $S$  in  $R^{n+t}$ . Let  $\pi$  be the projection  $(\mathbf{x}, \mathbf{y}) \rightarrow \mathbf{y}$ . By Theorem 4.2.4, the image  $\pi(S)$  of  $S$  by  $\pi$  is a semi-algebraic set. This rises an algorithmic question, known as one-block quantifier elimination, that aims to compute a semi-algebraic formula  $\Theta$  defining  $\pi(S)$ , i.e.,

$$\exists \mathbf{x} : \Phi(\mathbf{x}, \mathbf{y}) \text{ is true} \Leftrightarrow \Theta(\mathbf{y}) \text{ is true.}$$

This problem is one of the main topics of this thesis. We will go back to this problem and present our contributions in Chapter 6 of the thesis.

We now recall the topology over  $R^n$ . For any real closed field  $R$ , one can define a topology over  $R^n$  similar to the Euclidean topology over  $\mathbb{R}^n$ .

**Definition 4.2.5.** *Let  $\eta = (\eta_1, \dots, \eta_n) \in R^n$  and  $r \in R$ ,  $r > 0$ . The open ball  $B(\eta, r)$  centered at  $\eta$  of radius  $r$  is defined as the set*

$$\{(x_1, \dots, x_n) \mid \sqrt{(\eta_1 - x_1)^2 + \dots + (\eta_n - x_n)^2} < r\}.$$

*A subset  $S$  of  $R^n$  is called an open set if for any  $\eta \in S$ , there exists an open ball  $B(\eta, r)$  contained in  $S$ . A set is closed if its complement in  $R^n$  is open.*

*A map  $\varphi : S \rightarrow S'$  is continuous if the pre-image of any open subset of  $S'$  is an open subset in  $S$ .*

*Further, we refer to this topology over  $R^n$  as the Euclidean topology.*

Next, we define the notions of continuous semi-algebraic maps and semi-algebraic homeomorphism for the Euclidean topology over  $R^n$ .

**Definition 4.2.6.** *A semi-algebraic map  $\varphi : S_1 \rightarrow S_2$  is continuous if the pre-image of any open set of  $S_2$  is an open set of  $S_1$ .*

*Two semi-algebraic sets  $S_1$  and  $S_2$  are semi-algebraically homeomorphic if there exists a bijective continuous semi-algebraic map  $\varphi : S_1 \rightarrow S_2$  such that  $\varphi^{-1}$  is also a continuous semi-algebraic map.*

The theorem below, known as Hardt's triviality theorem, implies that the fibers of a given continuous semi-algebraic map can be classified into finitely many types.

**Theorem 4.2.7** (Hardt's triviality theorem, [97]). *Let  $S \subset R^n$ ,  $S' \subset R^p$  be semi-algebraic sets and  $\varphi : S \rightarrow S'$  be a continuous semi-algebraic map. Then there exists a finite partition of  $S'$  into semi-algebraic sets*

$$S' = \bigcup_{i=1}^s S'_i$$

so that for  $1 \leq i \leq s$  and any  $\eta_i \in S'_i$ ,  $S'_i \times \varphi^{-1}(\eta_i)$  is semi-algebraically homeomorphic to  $\varphi^{-1}(S'_i)$ .

Now we discuss two important properties of a topological space: connectedness and compactness. The classical definitions in  $\mathbb{R}$  cannot be ported completely to a general closed field  $R$ . Below we present the necessary revisions of these notions in the semi-algebraic context.

A subset of a topological space is connected if it cannot be written as the disjoint union of two open subsets. However, this definition leads to some difficulties while working over real closed fields as we observe in the example below.

**Example 4.2.8.** *The field  $\mathbb{R}_{\text{alg}}$  of real algebraic numbers is disconnected as it is the union of two open subsets  $(-\infty, \pi) \cap \mathbb{R}_{\text{alg}}$  and  $(\pi, \infty) \cap \mathbb{R}_{\text{alg}}$ .*

Therefore, the following notion of semi-algebraic connectedness is more useful.

**Definition 4.2.9.** *A semi-algebraic set  $S \subset R^n$  is semi-algebraically connected if  $S$  is not the disjoint union of two non-empty semi-algebraic sets that are both open in  $S$ .*

*A semi-algebraic set  $S$  is semi-algebraically path connected when for every  $a, b \in S$ , there exists a continuous semi-algebraic function  $\phi : [0, 1] \rightarrow S$  such that  $\phi(0) = a$  and  $\phi(1) = b$ .*

In Example 4.2.8, the subset  $(-\infty, \pi) \cap \mathbb{R}_{\text{alg}}$  is not a semi-algebraic set since  $\pi$  is transcendental over  $\mathbb{R}$ .

**Proposition 4.2.10** ([9, Prop. 3.9]). *A real closed field is semi-algebraically connected.*

Note that even in  $\mathbb{R}$ , the connectedness does not imply the path connectedness. By contrast, the notions of semi-algebraic connectedness and semi-algebraic path connectedness are equivalent.

**Proposition 4.2.11** ([9, Theorem 5.23]). *A semi-algebraic set  $S \subset R^n$  is semi-algebraically connected if and only if it is semi-algebraically path connected.*

**Proposition 4.2.12** ([9, Theorem 5.21]). *Any semi-algebraic set  $S \subset R^n$  has finitely many semi-algebraic connected components.*

When  $R = \mathbb{R}$ , the notions of semi-algebraic connectedness and connectedness coincide.

**Proposition 4.2.13** ([9, Theorem 5.22]). *A semi-algebraic set of  $\mathbb{R}^n$  is semi-algebraically connected if and only if it is connected in  $\mathbb{R}^n$ .*

Unlike  $\mathbb{R}$ , a closed and bounded semi-algebraic set in  $R$  is not necessarily compact. In this thesis, we replace, when needed, the notion of compactness for semi-algebraic sets by closed and bounded semi-algebraic sets.

**Example 4.2.14.** *The interval  $[0, 1]$  is not compact in  $\mathbb{R}_{\text{alg}}$ . The family  $(]0, a[ \cup ]b, 1])$  for  $0 < a < \pi/4 < b < 1$ ,  $a, b \in \mathbb{R}_{\text{alg}}$  is an open cover of  $[0, 1]$  by semi-algebraic subsets of  $\mathbb{R}_{\text{alg}}$  and it is impossible to extract a finite cover from it.*

The smoothness of a semi-algebraic set can be defined through differential geometry.

**Definition 4.2.15** (Smoothness). *A semi-algebraic set  $S \subset \mathbb{R}^n$  is smooth at a point  $\eta$  if there exists an open subset  $U \subset \mathbb{R}^n$  such that  $S \cap U$  is diffeomorphic to  $\mathbb{R}^d$  for some integer  $d$  which is called the local dimension of  $S$  at  $\eta$ .*

*The semi-algebraic set  $S$  is smooth if  $S$  is smooth at every point  $\eta \in S$ . Given a smooth semi-algebraic set  $S$ ,  $S$  is equidimensional of dimension  $d$  if every point  $\eta \in S$  has dimension  $d$ .*

The definition of proper maps below is similar to Definition 2.9.5 in the context of algebraic sets.

**Definition 4.2.16** (Proper maps). *Let  $S \subset R^{n+p}$  be a semi-algebraic set and  $\varphi : V \rightarrow R^p$  be a semi-algebraic map. The map  $\varphi$  is proper at a point  $\eta \in R^p$  if there exists a neighborhood  $O$  of  $\eta$  such that  $\varphi^{-1}(\overline{O})$  is bounded, where  $\overline{O}$  denotes the closure of  $O$  for the Euclidean topology over  $R^p$ .*

*If  $\varphi$  is proper everywhere on its image, we say that the map  $\varphi$  is proper.*

The following important theorem allows to decompose the into part of invariant topology.

**Theorem 4.2.17** (Thom's isotopy lemma, [47]). *Let  $S \subset R^n$  be a semi-algebraic set and  $\varphi : S \rightarrow R^p$  be the projection onto the last  $p$  coordinates. We assume that*

- *$S$  is smooth and equidimensional;*
- *$S$  is locally closed;*
- *The projection  $\varphi$  is a proper map.*

*Let  $S' \subset R^p$  be a semi-algebraic set that does not contain any critical values of the restriction of  $\varphi$  to  $S$ . Then, for any  $\eta \in S'$ ,  $\varphi^{-1}(S')$  is diffeomorphic to  $\varphi^{-1}(\eta) \times S'$ .*

Theorem 4.2.17 provides a boundary of semi-algebraic sets in the destination space  $R^p$  over which the fibers of  $\varphi$  are topologically invariant. Using algorithmic tools from computer algebra (Jacobian criterion, elimination using Gröbner bases), one can compute defining systems of this boundary following its description given in the theorem.

### 4.3 Puiseux series

In real algebraic geometry, many algorithms require the smoothness on the input semi-algebraic sets to work correctly or more efficiently. When the semi-algebraic sets taken as input are singular, a commonly used technique is the deformation, which consists of basically two steps:

- Deform the input to obtain smooth semi-algebraic sets using some sufficiently small values;
- Taking the limit of these values to 0 to obtain the results on the original input.

This procedure of deformation is made rigorously by introducing infinitesimals and Puiseux series, which we recall the definitions and properties in what follows.

We consider an infinitesimal  $\varepsilon$ , i.e., a transcendental element over  $\mathbb{R}$  such that  $0 < \varepsilon < r$  for any positive element  $r \in \mathbb{R}$ . The fields of Puiseux series over  $\mathbb{R}$  and  $\mathbb{C}$  are defined as follows.

**Definition 4.3.1.** *The field of Puiseux series over  $\mathbb{R}$ , denoted by  $\mathbb{R}\langle\varepsilon\rangle$ , is*

$$\mathbb{R}\langle\varepsilon\rangle = \left\{ \sum_{i \geq i_0} a_i \varepsilon^{i/q} \mid i \in \mathbb{N}, i_0 \in \mathbb{Z}, q \in \mathbb{N} - \{0\}, a_i \in \mathbb{R} \right\}.$$

Similarly, one defines  $\mathbb{C}\langle\varepsilon\rangle$  as for  $\mathbb{R}\langle\varepsilon\rangle$  but taking the coefficients of the series in  $\mathbb{C}$ .

We have that  $\mathbb{C}\langle\varepsilon\rangle = \mathbb{R}\langle\varepsilon\rangle[u] / \langle u^2 + 1 \rangle$ .

**Theorem 4.3.2** ([9, Theorem 2.113]). *The field  $\mathbb{R}\langle\varepsilon\rangle$  is a real closed field. As a consequence,  $\mathbb{C}\langle\varepsilon\rangle$  is an algebraic closure of  $\mathbb{R}\langle\varepsilon\rangle$ .*

**Definition 4.3.3.** *Given a Puiseux series*

$$\sigma = \sum_{i \geq i_0} a_i \varepsilon^{i/q} \in \mathbb{R}\langle\varepsilon\rangle$$

with  $a_{i_0} \neq 0$ ,  $a_{i_0}$  is called the initial coefficient of  $\sigma$ .

When  $i_0 \geq 0$ ,  $\sigma$  is said to be bounded over  $\mathbb{R}$ . The subset of  $\mathbb{R}\langle\varepsilon\rangle$  of elements which are bounded over  $\mathbb{R}$  is denoted by  $\mathbb{R}\langle\varepsilon\rangle_b$ .

The limit of a Puiseux series is defined algebraically by sending  $\varepsilon$  to 0.

**Definition 4.3.4.** *Let  $\lim_\varepsilon : \mathbb{R}\langle\varepsilon\rangle_b \rightarrow \mathbb{R}$  be the function that maps  $\sigma$  to  $a_0$  (which is 0 when  $i_0 > 0$ ) and writes  $\lim_\varepsilon \sigma = a_0$ . Note that  $\lim_\varepsilon$  is a ring homomorphism from  $\mathbb{R}\langle\varepsilon\rangle_b$  to  $\mathbb{R}$ .*

All the definitions above extend to  $\mathbb{R}\langle\varepsilon\rangle^n$  componentwise.

For a semi-algebraic set  $\mathcal{S} \subset \mathbb{R}\langle\varepsilon\rangle^n$ , we naturally define the limit of  $\mathcal{S}$  as

$$\lim_\varepsilon \mathcal{S} = \left\{ \lim_\varepsilon \mathbf{x} \mid \mathbf{x} \in \mathcal{S} \text{ and } \mathbf{x} \text{ is bounded over } \mathbb{R} \right\}.$$

For each semi-algebraic set defined over  $\mathbb{R}$ , we can associate to it a semi-algebraic in  $\mathbb{R}\langle\varepsilon\rangle^n$  using the notion of extension below.

**Definition 4.3.5.** Let  $S \subset \mathbb{R}^n$  be a semi-algebraic set defined by a semi-algebraic formula  $\Phi$  whose atoms lie in  $\mathbb{R}[x_1, \dots, x_n]$ . We denote by  $\text{ext}(S, \mathbb{R}\langle\varepsilon\rangle)$  the semi-algebraic set of points which are solutions of  $\Phi$  in  $\mathbb{R}\langle\varepsilon\rangle^n$ , i.e.,

$$\text{ext}(S, \mathbb{R}\langle\varepsilon\rangle) = \{\eta \in \mathbb{R}\langle\varepsilon\rangle^n \mid \Phi(\eta) \text{ is true}\}.$$

Note that the definition of  $\text{ext}(S, \mathbb{R}\langle\varepsilon\rangle)$  depends only on the semi-algebraic set  $S$  and not on the choice of the defining formula  $\Phi$  (see [9, Prop. 2.105]). We also have the extension of a semi-algebraic map.

**Definition-Proposition 4.3.6** ([9, Prop. 2.108]). Let  $S \subset \mathbb{R}^n$  and  $S' \subset \mathbb{R}^p$  be two semi-algebraic sets. Given a semi-algebraic map  $f : S \rightarrow S'$  whose graph is a semi-algebraic set  $G \subset S \times S'$ . The extension of  $f$  to  $\mathbb{R}\langle\varepsilon\rangle$ , denoted by  $\text{ext}(f, \mathbb{R}\langle\varepsilon\rangle)$ , is defined as a semi-algebraic map from  $\text{ext}(S, \mathbb{R}\langle\varepsilon\rangle)$  to  $\text{ext}(S', \mathbb{R}\langle\varepsilon\rangle)$  whose graph is  $\text{ext}(G, \mathbb{R}\langle\varepsilon\rangle)$ .

The following result is the foundation of the deformation technique we will use further in the thesis.

**Proposition 4.3.7** ([167, Lemma 3.5]). Let  $f \in \mathbb{C}[x_1, \dots, x_n]$  and  $\varepsilon$  be an infinitesimal. The algebraic set of  $\mathbb{C}\langle\varepsilon\rangle^n$  defined by  $f = \varepsilon$  (or  $f = -\varepsilon$ ) is a smooth algebraic set.

The following propositions will be useful for taking the limits of deformed semi-algebraic sets.

**Proposition 4.3.8** ([9, Prop. 12.49]). If  $S' \subset \mathbb{R}\langle\varepsilon\rangle^n$  is a semi-algebraic set, then  $\lim_\varepsilon S'$  is a closed semi-algebraic set. Moreover, if  $S' \subset \mathbb{R}\langle\varepsilon\rangle^n$  is a semi-algebraic set bounded over  $\mathbb{R}$  and semi-algebraically connected, then  $\lim_\varepsilon S'$  is semi-algebraically connected.

**Proposition 4.3.9** ([9, Prop. 12.51]). Given a polynomial  $f \in \mathbb{R}[x_1, \dots, x_n]$ , we introduce the following notations:

- $V$  denotes the real algebraic subset of  $\mathbb{R}^n$  defined by  $f = 0$ ;
- $V_\varepsilon$  denotes the real algebraic subset of  $\mathbb{R}\langle\varepsilon\rangle^n$  defined by  $f^2 = \varepsilon^2$ ;
- $V_b \subset \mathbb{R}^n$  denotes the union of the semi-algebraically connected components of  $V$  which are bounded over  $\mathbb{R}$ ;
- $V_{\varepsilon,b} \subset \mathbb{R}\langle\varepsilon\rangle^n$  denotes the union of the semi-algebraically connected components of  $V_\varepsilon$  which are bounded over  $\mathbb{R}$ .

Then we have

$$\lim_\varepsilon V_{\varepsilon,b} = V_b.$$

## 4.4 Real root counting

The fundamental theorem of algebra implies that any polynomial of degree  $D$  defined over an algebraically closed field  $K$  has exactly  $D$  solutions in  $K$  counted with multiplicities. However, for polynomials defined over a real closed field  $R$ , the possible number of solutions in  $R$  can vary between 0 and  $D$ . In this section, we recall some classical algorithmic tools for counting the number of real solutions of polynomial equations. Further in Chapter 5, we will extend some of these tools to the case of polynomial systems depending on parameters.

Throughout this section,  $R$  is a real closed field and  $K = R[T]/\langle T^2 + 1 \rangle$ .

### 4.4.1 Notations

We start with some definitions that will be commonly used throughout the section.

**Definition 4.4.1** (Sign variations). *The number of sign variations in a sequence  $a = (a_1, \dots, a_s)$  of elements in  $R \setminus \{0\}$  is defined by*

$$\text{Var}(a_1) = 0, \text{Var}(a_1, \dots, a_{i+1}) = \begin{cases} \text{Var}(a_1, \dots, a_i) + 1 & \text{if } a_i a_{i+1} < 0, \\ \text{Var}(a_1, \dots, a_i) & \text{if } a_i a_{i+1} > 0. \end{cases}$$

*The sign variations of a sequence containing zeros is defined as the sign variations of the same sequence with all zeros removed.*

*Let  $\mathbf{f} = f_1, \dots, f_s$  be a sequence of polynomials in  $R[u]$  and let  $a \in R \cup \{-\infty, \infty\}$ . The number of sign variations of  $\mathbf{f}$  at  $a$ , denoted by  $\text{Var}(\mathbf{f}; a)$  is  $\text{Var}(f_1(a), \dots, f_s(a))$  (at  $-\infty$  and  $\infty$ , the signs to consider are the signs of the leading terms of  $f_1, \dots, f_s$ ).*

*We define  $\text{Var}(p; a, b) = \text{Var}(p; a) - \text{Var}(p; b)$ .*

**Example 4.4.2.** *We have that  $\text{Var}(1, -2, 0, 0, 3, 4, 0, -5, 6) = \text{Var}(1, -2, 3, 4, -5, 6) = 4$ .*

**Definition 4.4.3** (Tarski's query). *Let  $\mathbf{f} = (f_1, \dots, f_s) \subset R[x_1, \dots, x_n]$  and  $g \in R[x_1, \dots, x_n]$ . Assume that the system of equations*

$$f_1 = \dots = f_s = 0$$

*has finitely many solutions in  $R^n$ . The Tarski's query of  $g$  for  $\mathbf{f}$  is defined as*

$$\begin{aligned} \text{TarskiQuery}(\mathbf{f}, g) &= \sum_{\eta \in V_R(\mathbf{f})} \text{sign } g(\eta) \\ &= |\{\eta \in V_R(\mathbf{f}) \mid g(\eta) > 0\}| - |\{\eta \in V_R(\mathbf{f}) \mid g(\eta) < 0\}|. \end{aligned}$$

### 4.4.2 Sturm sequences and Sturm-Habicht sequences

In this subsection, we recall the classical results of real root counting using Sturm sequences and Sturm-Habicht sequences.

Sturm's theorems allows us to count the number of real solutions of one polynomial  $p$  over an interval using the signed remainder sequence of  $p$  and its derivative  $p'$  defined below.

**Definition 4.4.4** (Signed remainder sequence). Given  $p, q \in R[u]$  with  $q \neq 0$ , we denote the remainder of the division of  $p$  to  $q$  by  $\text{Rem}(p, q)$ .

The signed remainder sequence  $\text{sRem}(p, q)$  of  $p$  and  $q$  is defined as follows:

- If  $q = 0$ , then  $\text{sRem}(p, q) = p$ .
- If  $p = 0$ , then  $\text{sRem}(p, q) = 0, q$ .
- If  $p, q \neq 0$  then

$$\text{sRem}(p, q) = r_0, \dots, r_s,$$

where  $r_0 = p, r_1 = q$  and for  $i \geq 1$ , if  $r_i$  does not divide  $r_{i-1}$  then

$$r_{i+1} = -\text{Rem}(r_{i-1}, r_i)$$

where  $s$  is defined by  $r_s \neq 0$  and  $\text{Rem}(r_{s-1}, r_s) = 0$ .

**Theorem 4.4.5** (Sturm's theorem, [188], [9, Theorem 2.62]). Let  $p \in K[u]$ . Given  $a$  and  $b$  in  $R \cup \{-\infty, +\infty\}$  that are not roots of  $p$ , we have that

$$\text{Var}(\text{sRem}(p, p'); a, b) = |\{\eta \in (a, b) \mid p(\eta) = 0\}|.$$

**Example 4.4.6.** Let  $p = u^4 - 5u^2 + 4$ . The signed remainder sequence of  $p$  and  $p'$  is

$$\begin{aligned} \text{sRem}_0(p, p') &= u^4 - 5u^2 + 4, \\ \text{sRem}_1(p, p') &= 4u^3 - 10u, \\ \text{sRem}_2(p, p') &= \frac{5}{2}u^2 - 4, \\ \text{sRem}_3(p, p') &= \frac{18}{5}u, \\ \text{sRem}_4(p, p') &= 4. \end{aligned}$$

The signs of the leading coefficients of the sequence above at  $+\infty$  and  $-\infty$  are respectively

$$(+, +, +, +, +) \quad \text{and} \quad (+, -, +, -, +).$$

Hence,  $\text{Var}(\text{sRem}(p, p'); -\infty, +\infty) = 4 - 0 = 4$ . The polynomial  $p$  has actually 4 real roots: 1, -1, 2, and -2.

A more general variant, known as Tarski's theorem, counts the number of real solutions of one polynomial  $p$  with respect to the sign condition of a polynomial  $q$  using the Sturm sequence of  $p$  and  $q$ , which we define below.

**Definition 4.4.7** (Sturm sequence). Let  $p, q \in R[u]$ . The Sturm sequence of  $p$  and  $q$ , denoted by  $\text{Sturm}(p, q)$ , is the signed remainder sequence  $\text{sRem}(p, p'q)$  where  $p'$  denotes the derivative of  $p$ .

**Theorem 4.4.8** ([9, Theorem 2.73]). *Let  $p, q \in K[u]$ . Given  $a$  and  $b$  in  $R \cup \{-\infty, +\infty\}$  that are not roots of  $p$ . Then, we have that*

$$\text{Var}(\text{Sturm}(p, q); a, b) = |\{\eta \in (a, b) \mid p(\eta) = 0, q(\eta) > 0\}| - |\{\eta \in (a, b) \mid p(\eta) = 0, q(\eta) < 0\}|.$$

Note that Sturm's theorem is a particular case of the theorem above by taking  $q = 1$ .

**Example 4.4.9.** *Taking  $p = u^4 - 5u^2 + 4$  (as in Example 4.4.6) and  $q = u^2 - 2$ , the Sturm sequence of  $p$  and  $q$ , defined as the signed remainder sequence of  $p$  and  $p' \cdot q$ , is*

$$\begin{aligned} \text{sRem}_0(p, p' \cdot q) &= u^4 - 5u^2 + 4, \\ \text{sRem}_1(p, p' \cdot q) &= 4u^5 - 18u^3 + 20u, \\ \text{sRem}_2(p, p' \cdot q) &= -u^4 + 5u^2 - 4, \\ \text{sRem}_3(p, p' \cdot q) &= -2u^3 - 4u, \\ \text{sRem}_4(p, p' \cdot q) &= -7u^2 + 4, \\ \text{sRem}_5(p, p' \cdot q) &= \frac{36}{7}u, \\ \text{sRem}_6(p, p' \cdot q) &= -4. \end{aligned}$$

*The signs of at  $+\infty$  and  $-\infty$  are respectively*

$$(+, +, -, -, -, +, -) \quad \text{and} \quad (+, -, -, +, -, -, -).$$

*Thus,  $\text{Var}(\text{Sturm}(p, q); a, b) = 3 - 3 = 0$ . This corresponds to*

$$\{\eta \in (a, b) \mid p(\eta) = 0, q(\eta) > 0\} = \{2, -2\} \quad \text{and} \quad \{\eta \in (a, b) \mid p(\eta) = 0, q(\eta) < 0\} = \{1, -1\}.$$

From the algorithmic aspect, using the signed remainder sequences (therefore, Sturm sequences) has several disadvantages. Firstly, the bit-size growth of signed remainder sequences is not well-controlled, which makes estimating the complexity of Sturm sequences difficult. We can observe this in the example below.

**Example 4.4.10.** *Given  $p = 75u^6 - 92u^5 + 6u^3 + 74u^2 + 72u + 37$  and  $q = -23u^4 + 8u^3 + 44u^2 + 29u + 98$ , the leading coefficients of  $\text{sRem}(p, q)$  is*

$$\begin{aligned} &\frac{-45660506}{12167}, \\ &\frac{11894666533333043}{2084881808176036}, \\ &\frac{-1424856298844988047595604912488}{11628428695585856081319204047}, \\ &\frac{-407687090204585473313605266765585776122109512697}{108197738617155415021320079773703340240647056}. \end{aligned}$$

Assume now that the coefficients of  $p$  and  $q$  depends on some parameters, the second weakness of signed remainder sequence is its lack of specialization properties. As computing the sequence  $\text{sRem}(p, q)$  depends on a division process that may involve dividing by some polynomials of parameters, one cannot specialize the signed remainder sequence at the values of parameters which cancel those divisors. Hence, one may need to repeat the whole computation to obtain the signed remainder sequences of specializations of  $p$  and  $q$ .

To overcome this inconvenience, in [88], the authors introduced Sturm-Habicht sequences, a variant of Sturm sequences which provide similar tools for counting the number of real solutions. These sequences, constructed through the signed subresultant coefficient sequence, are therefore determinant polynomials of certain matrices. As a consequence, they inherit a specialization property and well-controlled bit-size growth of the signed subresultant coefficient sequences.

**Definition 4.4.11** (Sylvester-Habicht matrix). *Let  $p, q \in R[u]$  of degree  $k$  and  $\ell$  respectively with  $k \geq \ell$ ,*

$$\begin{aligned} p &= a_k u^k + \cdots + a_0, \\ q &= b_\ell u^\ell + \cdots + b_0. \end{aligned}$$

*For  $0 \leq j \leq \ell$ , the  $j$ -th Sylvester-Habicht matrix of  $p$  and  $q$ , denoted by  $\text{SylHa}_j(p, q)$  is the matrix*

$$\begin{bmatrix} a_k & \dots & \dots & \dots & \dots & a_0 & 0 & 0 \\ 0 & \ddots & & & & & \ddots & 0 \\ \vdots & \ddots & a_k & \dots & \dots & \dots & \dots & a_0 \\ \vdots & & 0 & b_\ell & \dots & \dots & \dots & b_0 \\ \vdots & \ddots & & & & & & 0 \\ 0 & \ddots & & & & \ddots & \ddots & \vdots \\ b_\ell & \dots & \dots & \dots & b_0 & 0 & \dots & 0 \end{bmatrix}.$$

*It has  $k + \ell - j$  columns and  $k + \ell - 2j$  rows.*

*The signed subresultant coefficient sequence of  $p$  and  $q$ , denoted by  $\text{sRes}(p, q)$ , is the sequence*

$$\text{sRes}(p, q) = \text{sRes}_k(p, q), \dots, \text{sRes}_0(p, q),$$

*where*

- $\text{sRes}_k(p, q) = a_k$ ,
- $\text{sRes}_{k-1}(p, q) = b_\ell$ ,
- $\text{sRes}_j(p, q) = 0$  for  $\ell < j < k - 1$ ,

- For  $0 \leq j \leq \ell$ ,  $\text{sRes}_j(p, q)$  is the determinant of the square matrix obtained by taking the first  $k + \ell - 2j$  columns of  $\text{SylHa}_j(p, q)$ .

The signed subresultant coefficients can be computed through a variant of Euclidean algorithm of  $p$  and  $q$  in which the dividend of each step is multiplied by some coefficient to avoid introducing denominators in the division (see [9, Algo. 8.77]). Hence, these signed subresultant coefficients lie in the ring generated by the coefficients of the polynomials  $p$  and  $q$  and do not involve any denominator.

The following proposition gives a specialization property for signed subresultant sequences.

**Proposition 4.4.12** ([9, Prop. 8.74]). *Let  $\mathbf{y} = (y_1, \dots, y_t)$  be the parameters and  $\phi$  be the ring morphism from  $R[\mathbf{y}][u]$  to  $R[u]$  by assigning  $\mathbf{y}$  to a value  $\eta \in R^t$ . Given  $p, q \in R[\mathbf{y}][u]$  such that  $\deg(\phi(p)) = \deg(p)$  and  $\deg(\phi(q)) = \deg(q)$ , then for all  $j \leq p$ ,*

$$\text{sRes}_j(\phi(p), \phi(q)) = \phi(\text{sRes}_j(p, q)).$$

**Example 4.4.13.** *We consider for example the polynomial*

$$p = u^4 + y_1 u^2 + y_2 u + y_3.$$

*The signed subresultant sequence of  $p$  and  $p'$  is formed by the polynomials*

$$\begin{aligned} \text{sRes}_4(p, p') &= u^4 + y_1 u^2 + y_2 u + y_3, \\ \text{sRes}_3(p, p') &= 4u^3 + 2y_1 u + y_2, \\ \text{sRes}_2(p, p') &= -4(2y_1 u^2 + 3y_2 u + 4y_3), \\ \text{sRes}_1(p, p') &= 4((8y_1 y_3 - 9y_2^2 - 2y_1^3)u - y_1^2 y_2 - 12y_2 y_3), \\ \text{sRes}_0(p, p') &= 256y_3^3 - 128y_1^2 y_2^2 + 144y_1 y_2^2 y_3 + 16y_1^4 y_3 - 27y_2^4 - 4y_1^3 y_2^2. \end{aligned}$$

*When specialize  $a = 0$ , the subresultant sequence of  $p_0 = u^4 + y_2 u + y_3$  and  $p'_0 = 4u^3 + y_2$  is*

$$\begin{aligned} \text{sRes}_4(p_0, p'_0) &= u^4 + y_2 u + y_3, \\ \text{sRes}_3(p_0, p'_0) &= 4u^3 + y_2, \\ \text{sRes}_2(p_0, p'_0) &= -4(3y_2 u + 4y_3), \\ \text{sRes}_1(p_0, p'_0) &= 4(-9y_2^2 u - 12y_2 y_3), \\ \text{sRes}_0(p_0, p'_0) &= 256y_3^3 - 27y_2^4, \end{aligned}$$

*which agrees with the specialization of the sequence above at  $a = 0$ .*

Example 4.4.14 below illustrates the growth of bit-sizes in signed subresultant coefficients.

**Example 4.4.14.** We take the polynomials

$$p = 75u^6 - 92u^5 + 6u^3 + 74u^2 + 72u + 37 \quad \text{and} \quad q = -23u^4 + 8u^3 + 44u^2 + 29u + 98$$

in Example 4.4.10. The signed subresultant coefficients of  $p$  and  $q$  is

$$-529, -199561, 82908687856, 1542839439026884, -51592162747958902864,$$

which have smaller bit-sizes than the coefficients of the signed remainder sequence in Example 4.4.10.

Using the signed subresultant sequences, we define Sturm-Habicht sequence, a variant of Sturm sequence with better algorithmic properties.

**Definition 4.4.15** (Sturm-Habicht sequence). Let  $p$  and  $q$  be two polynomials in  $R[u]$  and  $\overline{p' \cdot q}$  be the remainder of  $p' \cdot q$  by  $p$ . The Sturm-Habicht sequence of  $p$  and  $q$  is the signed subresultant coefficient sequence  $\text{sRes}(p, \overline{p' \cdot q})$ .

To state the root counting theorem for Sturm-Habicht sequences, we need the following definition.

**Definition 4.4.16** (Permanences minus variations). Let  $a = a_0, \dots, a_s$  be a sequence of elements in  $R$  such that  $a_0 \neq 0$ . Let  $\ell < k$  such that  $a_{k-1} = \dots = a_{\ell+1} = 0$  and  $a_\ell \neq 0$  and  $\tilde{a}$  denotes the subsequence  $a_\ell, \dots, a_0$ .

$$\text{PmV}(a) = \begin{cases} 0 & \text{if } \tilde{a} = \emptyset, \\ \text{PmV}(\tilde{a}) + (-1)^{(k-\ell)(k-\ell-1)/2} \text{sign}(a_k a_\ell) & \text{if } k - \ell \text{ is odd,} \\ \text{PmV}(\tilde{a}) & \text{if } k - \ell \text{ is even,} \end{cases}$$

**Theorem 4.4.17** ([9, Theorem 4.32]). Let  $p$  and  $q$  be two polynomials in  $R[u]$  and  $\overline{p' \cdot q}$  be the remainder of  $p' \cdot q$  by  $p$ . Then

$$\text{PmV}(\text{sRes}(p, \overline{p' \cdot q})) = \text{TarskiQuery}(p, q; -\infty, +\infty).$$

**Example 4.4.18.** We continue with the polynomials

$$p = u^4 - 5u^2 + 4 \quad \text{and} \quad q = u^2 - 2.$$

The remainder of  $p' \cdot q$  by  $p$  is  $r = 2u^3 + 4u$ . Computing the signed subresultants of  $p$  and  $r$  gives

$$\begin{aligned} \text{sRes}_4(p, r) &= u^4 - 5u^2 + 4, \\ \text{sRes}_3(p, r) &= 2u^3 + 4u, \\ \text{sRes}_2(p, r) &= 28u^2 - 16, \\ \text{sRes}_1(p, r) &= 1008u, \\ \text{sRes}_0(p, r) &= 20736. \end{aligned}$$

We deduce that  $\text{PmV}(\text{sRes}(p, r)) = 0$ , which agrees with  $\text{TarskiQuery}(p, q; -\infty, +\infty) = 0$  as in Example 4.4.9.

### 4.4.3 Hermite quadratic forms

In [106], Hermite introduced a method for counting the solutions of a given univariate polynomial by associating to it a quadratic form. Later on, in [160], Hermite's quadratic forms were generalized to multivariate zero-dimensional systems.

Hermite's quadratic forms are the key ingredient for designing our algorithm for solving parametric polynomial systems in Chapter 5. In what follows, we recall the definition and basic properties of Hermite's quadratic forms.

Throughout this subsection,  $R$  is a real closed field and  $K = R[T]/\langle T^2 + 1 \rangle$ , which is algebraically closed.

Given a zero-dimensional ideal  $I \subset R[\mathbf{x}]$  where  $\mathbf{x} = (x_1, \dots, x_n)$ , the quotient ring  $R[\mathbf{x}]/I$  is a  $R$ -vector space of finite dimension (Theorem 3.3.1); its dimension is denoted as  $\delta$ .

We define the multiplication maps of  $R[\mathbf{x}]/I$  as follows.

**Definition 4.4.19.** For any  $p \in R[\mathbf{x}]$ , the multiplication map  $\mathcal{L}_p$  is defined as

$$\begin{aligned} \mathcal{L}_p : R[\mathbf{x}]/I &\rightarrow R[\mathbf{x}]/I, \\ \bar{q} &\mapsto \overline{p \cdot q}, \end{aligned}$$

where  $\bar{q}$  and  $\overline{p \cdot q}$  are respectively the classes of  $q$  and  $p \cdot q$  in the quotient ring  $R[\mathbf{x}]/I$ .

Note that the map  $\mathcal{L}_p$  is an endomorphism of  $R[\mathbf{x}]/I$  as a  $R$ -vector space.

We consider a basis  $B = \{b_1, \dots, b_\delta\}$  of  $R[\mathbf{x}]/I$ . Such a basis  $B$  can be derived from Gröbner bases as shown in Section 3.3. We fix an admissible monomial ordering  $\succ$  over the set of monomials in the variables  $\mathbf{x}$  and compute a Gröbner basis  $G$  with respect to the ordering  $\succ$  of the ideal  $I$ . Then, the monomials that are not divisible by any leading monomial of elements of  $G$  form a basis of  $R[\mathbf{x}]/I$ .

Recall that, for an element  $p \in R[\mathbf{x}]$ , we denote by  $\bar{p}$  the class of  $p$  in the quotient ring  $R[\mathbf{x}]/I$ . A representative of the class  $\bar{p}$  can be derived by computing the normal form of  $p$  by the Gröbner basis  $G$ , which results in a linear combination of elements of  $B$  with coefficients in  $R$ .

For any  $p \in R[\mathbf{x}]$ , the multiplication map  $\mathcal{L}_p$  is an endomorphism of  $R[\mathbf{x}]/I$ . Thus, it admits a matrix representation with respect to  $B$ , whose entries are elements in  $R$ .

**Example 4.4.20.** Let  $I = \langle x_1^2 + 2x_1x_2 + 3x_1 + x_2 + 1, x_2^2 + x_1x_2 + 2x_2 + 1 \rangle$ .

The reduced Gröbner basis of  $I$  with respect to the  $\text{grevlex}(x_1 \succ x_2)$  ordering is

$$\{x_2^3 + 2x_2^2 + x_1 + 2x_2 + 1, x_1^2 - 2x_2^2 + 3x_1 - 3x_2 - 1, x_1x_2 + x_2^2 + 2x_2 + 1\}.$$

Then, we derive a basis  $B = \{1, x_2, x_1, x_2^2\}$  for the  $R$ -vector space  $R[x_1, x_2]/I$ . The multiplication map  $\mathcal{L}_1$  is the identity endomorphism of  $R[x_1, x_2]/I$  and therefore it is represented by an identity matrix of size  $4 \times 4$ .

We compute the matrix representation of  $\mathcal{L}_{x_1}$  in  $B$  through these normal form reductions below:

$$\begin{aligned}x_1 \cdot 1 &= x_1, \\x_1 \cdot x_2 &= -1 - 2x_2 - x_2^2, \\x_1 \cdot x_1 &= 1 + 3x_2 - 3x_1 + 2x_2^2, \\x_1 \cdot x_2^2 &= 1 + x_2 + x_1.\end{aligned}$$

Thus, the matrix represents  $\mathcal{L}_{x_1}$  with respect to  $B$  is

$$\begin{bmatrix} 0 & 0 & 1 & 0 \\ -1 & -2 & 0 & -1 \\ 1 & 3 & -3 & 2 \\ 1 & 1 & 1 & 0 \end{bmatrix}.$$

Let  $V(I)$  be the finite set of zeros in  $K^n$  of  $I$ . For each  $\eta \in V(I)$ , we denote by  $\mu(\eta)$  the multiplicity of  $\eta$  defined in Proposition 3.3.3. The eigenvalues of these multiplication maps provide information on the zeros of  $I$ .

**Theorem 4.4.21** ([9, Theorem 4.98]). *Let  $p \in R[x_1, \dots, x_n]$ . The eigenvalues of the multiplication map  $\mathcal{L}_p$  are the  $p(\eta)$ 's for  $\eta \in V(I)$ , with the multiplicity  $\mu(\eta)$ .*

**Theorem 4.4.22** (Stickelberger's theorem, [9, Theorem 4.99]). *Let  $I$  be a zero-dimensional ideal of  $R[x_1, \dots, x_n]$  and  $p \in R[x_1, \dots, x_n]$ . The linear map  $\mathcal{L}_p$  of  $R[x_1, \dots, x_n]/I$  has the following properties:*

- The trace of  $\mathcal{L}_p$  is

$$\text{trace}(\mathcal{L}_p) = \sum_{\eta \in V(I)} \mu(\eta) \cdot p(\eta).$$

- The determinant of  $\mathcal{L}_p$  is

$$\det(\mathcal{L}_p) = \prod_{\eta \in V(I)} p(\eta)^{\mu(\eta)}.$$

- The characteristic polynomial  $\chi(I, p, T)$  of  $\mathcal{L}_p$  is

$$\chi(I, p, T) = \prod_{\eta \in V(I)} (T - p(\eta))^{\mu(\eta)}.$$

Now we define Hermite's quadratic forms as follows.

**Definition 4.4.23** (Hermite's quadratic forms). *Let  $g \in R[x_1, \dots, x_n]$ . The Hermite quadratic form associated to  $I$  and  $g$  is defined as the bilinear form from  $R[\mathbf{x}]/I \times R[\mathbf{x}]/I$  to  $R[\mathbf{x}]/I$  that sends*

$$(\bar{p}, \bar{q}) \mapsto \text{trace}(\mathcal{L}_{p \cdot q \cdot g}),$$

where  $\text{trace}(\mathcal{L}_{p \cdot q \cdot g})$  is the trace of  $\mathcal{L}_{p \cdot q \cdot g}$  as an endomorphism of  $R[\mathbf{x}]/I$ .

With a fixed basis  $B = \{b_1, \dots, b_\delta\}$  of the  $R$ -vector space  $R[\mathbf{x}]/I$ , Hermite's quadratic form of the ideal  $I$  admits a matrix representation with respect to  $B$ . Thus, we also have the definition of Hermite matrix of  $I$  with respect to the basis  $B$ .

**Definition 4.4.24** (Hermite matrix). *Let  $I \subset R[x_1, \dots, x_n]$  be a zero-dimensional ideal of degree  $\delta$  and  $g \in R[x_1, \dots, x_n]$ . For a basis  $B = \{b_1, \dots, b_\delta\}$  of  $R[\mathbf{x}]/I$  as  $R$ -vector space, we define the Hermite matrix of  $I$  with respect to the basis  $B$  as the symmetric matrix  $H = (h_{i,j})_{1 \leq i,j \leq \delta}$  where*

$$h_{i,j} = \text{trace}(\mathcal{L}_{b_i \cdot b_j \cdot g}).$$

**Example 4.4.25.** *When the ideal  $I$  is in shape position, using the lexicographic ordering,  $B$  is chosen to be  $\{1, x_n, \dots, x_n^{\delta-1}\}$ . Thus, the Hermite matrix  $H = (h_{i,j})_{1 \leq i,j \leq \delta}$  of  $I$  and any  $g \in R[x_1, \dots, x_n]$  with respect to this basis has the entry*

$$h_{i,j} = \text{trace} \left( \mathcal{L}_{x_n^{i+j-2} \cdot g} \right).$$

Note that  $H$  is a Hankel matrix.

**Example 4.4.26.** *Given the ideal  $I = \langle x_1^2 + x_2x_1 + 2x_2 + 3, x_2^2 + 2x_1x_2 + 3x_1 + 1 \rangle$ , the reduced Gröbner basis of  $I$  with respect to the  $\text{lex}(x_1 \succ x_2)$  ordering is*

$$G_{\text{lex}} = \{13x_1 + 2x_2^3 - 13x_2^2 - 46x_2 - 33, x_2^4 - 5x_2^3 - 36x_2^2 - 51x_2 - 28\}.$$

Hence, the basis of  $\mathbb{C}[x_1, x_2]/I$  associated to this Gröbner basis is  $B_{\text{lex}} = \{1, x_2, x_2^2, x_2^3\}$ . The Hermite matrix of  $I$  with respect to  $B_{\text{lex}}$  is

$$H_{\text{lex}} = \begin{bmatrix} 4 & 5 & 97 & 818 \\ 5 & 97 & 818 & 7949 \\ 97 & 818 & 7949 & 74280 \\ 818 & 7947 & 74280 & 701998 \end{bmatrix}.$$

On the other hand, computing the reduced Gröbner basis of the ideal  $I$  with respect to the  $\text{grevlex}(x_1 \succ x_2)$  ordering, we obtain

$$G_{\text{grevlex}} = \{2x_2^3 - 13x_2^2 + 13x_1 - 46x_2 - 33, 2x_1^2 - x_2^2 - 3x_1 + 4x_2 + 5, 2x_1x_2 + x_2^2 + 3x_1 + 1\}.$$

From this Gröbner basis, we derive the basis  $B_{\text{grevlex}} = \{1, x_2, x_1, x_2^2\}$  of  $\mathbb{C}[x_1, x_2]/I$  and, thus, the associated Hermite matrix is

$$H_{\text{grevlex}} = \begin{bmatrix} 4 & 5 & -1 & 97 \\ 5 & 97 & -49 & 818 \\ -1 & -49 & 27 & -338 \\ 97 & 818 & -338 & 7949 \end{bmatrix}.$$

We remark that the bit-sizes of entries of  $H_{\text{grevlex}}$  is smaller than the ones of  $H_{\text{lex}}$ . A possible explanation is that, the degree of polynomials in  $G_{\text{grevlex}}$  is slightly smaller than the one in  $G_{\text{lex}}$ , which leads to the difference in computing the entries through normal form reductions.

Finally, the theorem below shows how to use Hermite's matrix to count the number of solutions of a given zero-dimensional system.

**Proposition 4.4.27** ([9, Theorem 4.102]). *Let  $\mathbf{f} = (f_1, \dots, f_s) \in R[x_1, \dots, x_n]$  such that the ideal  $\langle f_1, \dots, f_s \rangle$  is zero-dimensional and  $g \in R[x_1, \dots, x_n]$ . Let  $H$  be the Hermite matrix of  $\mathbf{f}$  with respect to some basis of  $R[x_1, \dots, x_n]/\langle f_1, \dots, f_s \rangle$ . Then, we have*

- *The rank of  $H$  equals the number of distinct solutions in  $K^n$  of  $\mathbf{f}$  at which  $g$  is non-zero.*
- *The signature, i.e., the difference between the numbers of positive and negative eigenvalues, of  $H$  equals to the Tarski query  $\text{TarskiQuery}(\mathbf{f}, g)$ .*

When  $g$  is identically 1, we call  $H$  the Hermite matrix associated to  $\mathbf{f}$  for short. In this case, we have immediately that

- The rank of  $H$  equals the number of distinct solutions in  $K^n$  of  $\mathbf{f}$ .
- The signature of  $H$  equals to the number of distinct solutions in  $R^n$  of  $\mathbf{f}$ .

**Example 4.4.28.** *We continue with Example 4.4.26. The ideal  $I$  has 4 distinct complex roots and 2 distinct real roots. The rank and signature of both  $H_{\text{lex}}$  and  $H_{\text{grevlex}}$  are respectively 4 and 2.*

**Part II**

**Contributions**

# Chapter 5

## Real root classification algorithms

**Abstract.** In this chapter, we design a new algorithm for solving parametric systems of equations having finitely many complex solutions for generic values of the parameters.

More precisely, let  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  where  $\mathbf{x} = (x_1, \dots, x_n)$  are variables and  $\mathbf{y} = (y_1, \dots, y_t)$  are parameters. We denote by  $\mathcal{V} \subset \mathbb{C}^{n+t}$  the algebraic set defined by the simultaneous vanishing of the  $f_i$ 's and  $\pi$  is the projection

$$\pi : \mathbb{C}^{n+t} \rightarrow \mathbb{C}^t, \quad (\mathbf{x}, \mathbf{y}) \mapsto \mathbf{y}.$$

Under the assumptions that  $\mathbf{f}$  admits finitely many complex solutions when specializing  $\mathbf{y}$  to generic values and that the ideal generated by  $\mathbf{f}$  is radical, we solve the following algorithmic problem.

On input  $\mathbf{f}$ , we compute *semi-algebraic formulas* defining open semi-algebraic sets  $\mathcal{S}_1, \dots, \mathcal{S}_\ell$  in the parameters' space  $\mathbb{R}^t$  such that  $\cup_{i=1}^\ell \mathcal{S}_i$  is dense in  $\mathbb{R}^t$  and, for  $1 \leq i \leq \ell$ , the number of real points in  $\mathcal{V} \cap \pi^{-1}(\eta)$  is invariant when  $\eta$  ranges over  $\mathcal{S}_i$ .

Our algorithm exploits special properties of some well chosen monomial bases in the quotient algebra  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]/I$  where  $I \subset \mathbb{Q}(\mathbf{y})[\mathbf{x}]$  is the ideal generated by  $\mathbf{f}$  in  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$  as well as the specialization property of the so-called parametric Hermite matrices. This allows us to obtain “compact” representations of the semi-algebraic sets  $\mathcal{S}_i$ . These representations, encoded by minors of Hermite matrices, have a determinantal nature and are easy to evaluate.

When  $\mathbf{f}$  satisfies extra genericity assumptions (such as regularity), we use the theory of Gröbner bases to derive complexity bounds both on the number of arithmetic operations in  $\mathbb{Q}$  and the degree of the output polynomials. More precisely, letting  $D$  be the maximal degrees of the  $f_i$ 's and  $\mathfrak{D} = n(D-1)D^n$ , we prove that, on a generic input  $\mathbf{f} = (f_1, \dots, f_n)$ , one can compute those semi-algebraic formulas with

$$O \sim \left( \binom{t + \mathfrak{D}}{t} 8^t n^{2t+1} D^{3nt+2(n+t)+1} \right)$$

arithmetic operations in  $\mathbb{Q}$  and that the polynomials involved in these formulas have degree bounded by  $\mathfrak{D}$ .

Note that the state-of-the-art software for real root classification rely on algorithms which compute a CAD in the parameter space  $\mathbb{R}^t$  to obtain the semi-algebraic formulas. Even though there is no existing complexity analysis for these algorithms, the complexity of computing CAD alone is already doubly exponential in  $t$ . Hence, our algorithm has a better theoretical complexity of  $D^{O(nt)}$ .

We report on practical experiments which illustrate the efficiency of our algorithm, both on generic parametric systems and parametric systems coming from applications since it allows us to solve systems which are out of reach on the software of the state-of-the-art.

This is joint-work with M. Safey El Din.

## 5.1 Introduction

### 5.1.1 Problem statement

Let  $\mathbf{f} = (f_1, \dots, f_s)$  be a sequence of  $s$  polynomials in  $\mathbb{Q}[\mathbf{y}][\mathbf{x}]$  where the indeterminates  $\mathbf{y} = (y_1, \dots, y_t)$  are considered as *parameters* and  $\mathbf{x} = (x_1, \dots, x_n)$  are considered as *variables*. We denote by  $\mathcal{V} \subset \mathbb{C}^{n+t}$  the complex algebraic set defined by

$$f_1 = \dots = f_s = 0$$

and by  $\mathcal{V}_{\mathbb{R}}$  its real trace  $\mathcal{V} \cap \mathbb{R}^{n+t}$ . We consider also the projection on the parameter space  $\mathbf{y}$

$$\pi : \begin{array}{ccc} \mathbb{C}^n \times \mathbb{C}^t & \rightarrow & \mathbb{C}^t, \\ (\mathbf{x}, \mathbf{y}) & \mapsto & \mathbf{y}. \end{array}$$

Further, we say that  $\mathbf{f}$  satisfies Assumption (5.A) when the following holds.

**Assumption 5.A.** *There exists a non-empty Zariski open subset  $\mathcal{O} \subset \mathbb{C}^t$  such that  $\pi^{-1}(\eta) \cap \mathcal{V}$  is non-empty and finite for any  $\eta \in \mathcal{O}$ .*

In other words, assuming (5.A) ensures that, for a generic value  $\eta$  of the parameters, the sequence  $\mathbf{f}(\eta, \cdot)$  defines a finite algebraic set and hence finitely many real points. Note that, it is easy to prove that one can choose  $\mathcal{O}$  in a way that the number of complex solutions to the entries of  $\mathbf{f}(\eta, \cdot)$  is invariant when  $\eta$  ranges over  $\mathcal{O}$  (e.g. using the theory of Gröbner basis). This is no more the case when considering real solutions whose number may vary when  $\eta$  ranges over  $\mathcal{O}$ .

By Hardt's triviality theorem (Theorem 4.2.7), there exists a real algebraic *proper* subset  $\mathcal{R}$  of  $\mathbb{R}^t$  such that, for any non-empty connected open set  $\mathcal{U}$  of  $\mathbb{R}^t \setminus \mathcal{R}$  and  $\eta \in \mathcal{U}$ ,  $\pi^{-1}(\eta) \times \mathcal{U}$  is homeomorphic with  $\pi^{-1}(\mathcal{U})$ .

This leads us to consider the following real root classification problem.

**Problem RRC** (Real root classification). *On input  $\mathbf{f}$  satisfying Assumption (5.A), compute semi-algebraic formulas defining semi-algebraic sets  $\mathcal{S}_1, \dots, \mathcal{S}_\ell$  such that*

- (i) *The number of real points in  $\mathcal{V} \cap \pi^{-1}(\eta)$  is invariant when  $\eta$  ranges over  $\mathcal{S}_i$ , for  $1 \leq i \leq \ell$ ;*
- (ii) *The union of the  $\mathcal{S}_i$ 's is dense in  $\mathbb{R}^t$ ;*

*as well as at least one sample point  $\eta_i$  in each  $\mathcal{S}_i$  and the corresponding number of real points in  $\mathcal{V} \cap \pi^{-1}(\eta_i)$ .*

A collection of semi-algebraic formulas sets is said to solve Problem (RRC) for the input  $\mathbf{f}$  if it defines a collection of semi-algebraic sets  $\mathcal{S}_i$  satisfies the above properties (i) and (ii).

Our output will have the form  $\{(\Phi_i, \eta_i, r_i) \mid 1 \leq i \leq \ell\}$  where  $\Phi_i$  is a semi-algebraic formula defining the set  $\mathcal{S}_i$ ,  $\eta_i \in \mathbb{Q}^t$  is a sample point of  $\mathcal{S}_i$  and  $r_i$  is the corresponding number of real roots.

A weak version of Problem (RRC) would be to compute only a set  $\{\eta_1, \dots, \eta_\ell\}$  of sample points for a collection of semi-algebraic sets  $\mathcal{S}_i$  solving Problem (RRC) and their corresponding numbers of real points in  $\mathcal{V} \cap \pi^{-1}(\eta_j)$ .

**Example 5.1.1.** Consider the equation  $x^2 + y_1x + y_2 = 0$  where  $y_1$  and  $y_2$  are the parameters and  $x$  is the unique variable. While  $y_1^2 - 4y_2 \neq 0$ , this equation always has exactly two distinct complex solutions. On the other hand, its number of distinct real solutions can take any value from 0 to 2, depending on the sign of the discriminant  $y_1^2 - 4y_2$ . One possible output for Problem (RRC) on this toy example is the following:

$$\begin{cases} y_1^2 - 4y_2 < 0, & (0, 1) & : & 0 \text{ real solution,} \\ y_1^2 - 4y_2 = 0, & (2, 1) & : & 1 \text{ real solution,} \\ y_1^2 - 4y_2 > 0, & (1, 0) & : & 2 \text{ real solutions.} \end{cases}$$

Note that another possible output is

$$\begin{cases} y_1^2 - 4y_2 < 0, & (0, 1) & : & 0 \text{ real solution,} \\ y_1^2 - 4y_2 > 0, & (1, 0) & : & 2 \text{ real solutions.} \end{cases}$$

as the above two inequalities define semi-algebraic sets whose union is dense in  $\mathbb{R}^2$ . In Fig. 5.1, the green region is defined by  $y_1^2 - 4y_2 < 0$  and the red one corresponds to  $y_1^2 - 4y_2 > 0$ . We see that each fiber of  $\pi$  over the red region intersects the real surface defined by  $x^2 + y_1x + y_2 = 0$  at two distinct points. Whereas, the fibers over the green region are empty.

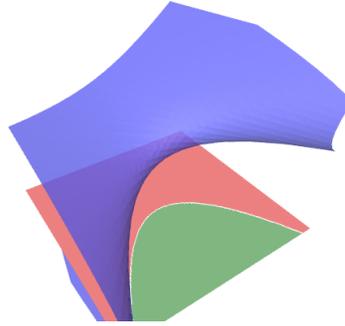


Figure 5.1: Real root classification of  $x^2 + y_1x + y_2 = 0$ .

A weak output consisting of only sample points is therefore

$$\begin{cases} (0, 1) & : & 0 \text{ real solution,} \\ (1, 0) & : & 2 \text{ real solutions.} \end{cases}$$

Problem (RRC) appears in many areas of engineering sciences such as robotics or medical imagery (see, e.g., [204, 45, 205, 64, 21]). In those applications, the behavior of mechanisms or complex systems depends on intrinsic parameters that are related by polynomial equations or inequalities. Thus, the polynomial systems arising from those applications are naturally parametric and most of the time the end-user is interested in classifying the number of real roots with respect to parameters' values.

### 5.1.2 Main results

As explained in Subsection 1.2.3, the state-of-the-art software for real root classification rely on cylindrical algebraic decomposition. Hence, their complexities are doubly exponential in  $t$ . In this chapter, we improve the state-of-the-art by designing algorithms of complexity  $D^{O(nt)}$  where  $D$  is a bound on the degree of input polynomials.

We start by revisiting methods based on Sturm-Habicht sequences in a multivariate context. We basically use the parametric geometric resolution of [178] to compute a rational parametrization of  $\mathcal{V} = V(\mathbf{f})$  with respect to the  $\mathbf{x}$ -variables. More precisely, we compute a sequence of polynomials  $(w, v_1, \dots, v_n)$  in  $\mathbb{Q}(\mathbf{y})[u]$  where  $u$  is a new variable, such that the constructible set  $\mathcal{Z} \subset \mathbb{C}^t \times \mathbb{C}^n$  of every point

$$\left( \eta, \frac{v_1}{\partial w / \partial u}(\eta, \vartheta), \dots, \frac{v_n}{\partial w / \partial u}(\eta, \vartheta) \right),$$

where  $(\eta, \vartheta) \in \mathbb{C}^t \times \mathbb{C}$  such that  $w(\eta, \vartheta) = 0$  and  $\eta$  does not cancel  $\partial w / \partial u$  and any denominator of  $(w, v_1, \dots, v_n)$ , is Zariski dense in  $\mathcal{V}$ , i.e., the Zariski closure of  $\mathcal{Z}$  coincides with  $\mathcal{V}$ .

Then, using the bi-rational equivalence between  $\mathcal{Z}$  and its projection on the  $(u, \mathbf{y})$ -space, we establish that semi-algebraic formulas solving Problem (RRC) can be obtained through the computation of the subresultant sequence associated to  $(w, \frac{\partial w}{\partial u})$ . This is admittedly *folklore* in symbolic computation but, as far as we know, is not explicitly written in the literature. In particular, the analysis of degree bounds derived from this strategy is one of our contributions.

Under some genericity assumptions on the input system, Theorem 5.1.2 establishes the complexity result of our Sturm-Habicht algorithm and also the degree bound for polynomials involved in the semi-algebraic formulas solving Problem (RRC) obtained this way. Its proof is given in Section 5.3, where all the genericity assumptions are clarified.

**Theorem 5.1.2.** *Let  $\mathbf{f} = (f_1, \dots, f_n) \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  be a generic parametric system and  $D$  be the largest total degree among the  $\deg(f_i)$ 's.*

*Then, there exists a probabilistic algorithm that computes semi-algebraic descriptions of a set of semi-algebraic sets solving Problem RRC within*

$$O \sim \left( \binom{t + 2D^{2n}}{t} 2^{5t} D^{5nt+3n} \right)$$

*arithmetic operations in  $\mathbb{Q}$  in case of success.*

These semi-algebraic formulas computed by this algorithm involve polynomials in  $\mathbb{Q}[\mathbf{y}]$  of degree bounded by  $2D^{2n}$ .

When reporting on experimental results, we will see that, even though the complexity bound we obtain lies in  $D^{O(nt)}$ , this approach does not allow us to solve problems faster than the state-of-the-art. One bottleneck comes from the fact that the polynomials of the output semi-algebraic formulas have degree way higher than the bound  $n(D - 1)D^n$  which we will prove to apply under the same assumptions as Theorem 5.1.2 using different algorithmic strategies.

Note that the above approach as well as the ones which compute polynomials in  $\mathbb{Q}[\mathbf{y}]$  to define boundaries of semi-algebraic sets in  $\mathbb{R}^t$  enjoying the properties needed to solve Problem (RRC) combine two steps of algebraic elimination. The semi-algebraic formulas are obtained through intermediate data who have been obtained through the first elimination step.

The rest of the chapter focuses on an alternative approach which computes semi-algebraic formulas solving Problem (RRC) by avoiding interlaced algebraic elimination steps. We will see (as announced earlier) that under genericity assumptions, this allows us to obtain a degree bound and an arithmetic cost which are better than the algorithm based on Sturm-Habicht sequences by one order of magnitude.

To do that, we rely on well-known properties of *Hermite quadratic forms* to count the real roots of zero-dimensional ideals (see Subsection 4.4.3).

Given a zero-dimensional ideal  $I \subset \mathbb{Q}[\mathbf{x}]$ , Hermite's quadratic form operates on the finite dimensional  $\mathbb{Q}$ -vector space  $A := \mathbb{Q}[\mathbf{x}]/I$  as follows:

$$\begin{aligned} A \times A &\rightarrow \mathbb{Q} \\ (h, k) &\mapsto \text{trace}(\mathcal{L}_{h,k}), \end{aligned}$$

where  $\mathcal{L}_{h,k}$  denotes the endomorphism  $p \mapsto h \cdot k \cdot p$  of  $A$ .

The number of distinct real (resp. complex) roots of the algebraic set defined by  $I$  equals the signature (resp. rank) of Hermite's quadratic form (Proposition 4.4.27). Recall that such a quadratic form is represented by a symmetric Hermite matrix of size  $\delta \times \delta$ , where  $\delta$  is the degree of  $I$ , once a basis of the finite dimensional vector space on which the form operates is fixed. Hence, the signature of a Hermite quadratic form can be computed from this matrix representation.

We first slightly extend the definition of Hermite's quadratic forms and Hermite's matrices to the context of parametric systems; we call them parametric Hermite quadratic forms and parametric Hermite matrices. This is easily done since the ideal of  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$  generated by  $\mathbf{f}$ , considering  $\mathbb{Q}(\mathbf{y})$  as the base field, has dimension zero. We also establish natural specialization properties for these parametric Hermite matrices.

Hence, a parametric Hermite matrix, similar to its zero-dimensional counterpart, allows one to count respectively the number of distinct real and complex roots at any parameters outside a proper algebraic sets of  $\mathbb{R}^t$  by evaluating the signature and rank of its specialization.

Based on this specialization property, we design two algorithms for solving Problem (RRC) and also its weak version for the input system  $\mathbf{f}$  which satisfies Assumption (5.A) and generates a radical ideal.

Our algorithm for the weak version of Problem (RRC) reduces to the following main steps.

- (a) Compute a parametric Hermite matrix  $\mathcal{H}$  associated to  $\mathbf{f} \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$ .
- (b) Compute a set of sample points  $\{\eta_1, \dots, \eta_\ell\}$  in the connected components of the semi-algebraic set of  $\mathbb{R}^t$  defined by  $\mathbf{w} \neq 0$  where  $\mathbf{w}$  is derived from  $\mathcal{H}$ .

This is done through the so-called critical point method (see e.g. [9, Chap. 12] and references therein) which are adapted to obtain practically fast algorithms following [171]. We will explain in detail this step in Section 5.2.

This algorithm takes as input  $m$  polynomials of degree  $D$  involving  $t$  variables and computes sample points per connected components in the semi-algebraic set defined by the non-vanishing of these polynomials using

$$O\left(\binom{D+t}{t} m^{t+1} 2^{3t} D^{2t+1}\right).$$

- (c) Compute the number  $r_i$  of real points in  $\mathcal{V} \cap \pi^{-1}(\eta_i)$  for  $1 \leq i \leq \ell$ .

This is done by simply evaluating the signature of the specialization of  $\mathcal{H}$  at each  $\eta_i$ .

It is worth noting that, in the algorithm above, we obtain through parametric Hermite matrices a polynomial  $\mathbf{w}$  that plays the same role as the discriminant varieties of [134] or the border polynomials of [202]. We will see in the section reporting experiments that our approach outperforms the other two for computing such a discriminating polynomial on every example we consider.

To return semi-algebraic formulas, we follow a slightly different routine:

- (a) Compute a parametric Hermite matrix  $\mathcal{H}$  associated to  $\mathbf{f} \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$ .
- (b) Compute a set of sample points  $\{\eta_1, \dots, \eta_\ell\}$  in the connected components of the semi-algebraic set of  $\mathbb{R}^t$  defined by  $\bigwedge_{i=1}^{\delta} M_i \neq 0$  where the  $M_i$ 's are the leading principal minors of  $\mathcal{H}$ . Again, this is done by the algorithm given in Section 5.2.
- (c) For  $1 \leq i \leq \ell$ , evaluate the sign pattern of  $(M_1, \dots, M_\delta)$  at the sample point  $\eta_i$ . From this sign pattern, we obtain a semi-algebraic formula representing the connected component corresponding to  $\eta_i$ .
- (d) Compute the number  $r_i$  of real points in  $\mathcal{V} \cap \pi^{-1}(\eta_i)$  for  $1 \leq i \leq \ell$ .

In Subsection 5.4.4, we will make clear how to perform Step (a) and present some remarks for optimization. For this, we rely on the theory of Gröbner bases and specialization properties similar to those already proven in [118]. This leaves some freedom when running the algorithm: since we rely on Gröbner bases, one may choose monomial orderings which are more convenient for practical computations.

In particular, the monomial basis of the quotient ring  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]/I$  where  $I$  is the ideal generated by  $\mathbf{f}$  in  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$  depends on the choice of the monomial ordering used for Gröbner bases computations. We describe the behavior of our algorithm when choosing the graded reverse lexicographical ordering whose interest for practical computations is explained in [13]. Further, we denote by  $\text{grevlex}(\mathbf{x})$  the graded reverse lexicographical ordering applied to the sequence of the variables  $\mathbf{x} = (x_1, \dots, x_n)$  (with  $x_1 \succ \dots \succ x_n$ ). Further, we also denote by  $\text{lex}(\mathbf{x})$  the lexicographical ordering  $x_1 \succ \dots \succ x_n$ .

We report, at the end of the chapter, on the practical behavior of this algorithm. In particular, it allows us to solve instances of Problem (RRC) which were not tractable by the state-of-the-art as well as the actual degrees of the polynomials in the output formula which are bounded by  $n(D-1)D^n$ . Using this algorithm, we successfully solve the application of Kuramoto model for 4 oscillators. As far as we know, this is the first symbolic solution for this application, comparing to [99] in which a numerical solution is given.

We actually prove such a statement under some genericity assumptions. Our main complexity result is stated below. Its proof is given in Subsection 5.6.2, where the genericity assumptions in use are given explicitly.

**Theorem 5.1.3.** *Let  $\mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}$  be the set of polynomials in  $\mathbb{C}[\mathbf{x}, \mathbf{y}]$  having total degree bounded by  $D$  and set  $\mathfrak{D} = n(D-1)D^n$ .*

*There exists a non-empty Zariski open set  $\mathcal{F} \subset \mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^n$  such that for  $\mathbf{f} = (f_1, \dots, f_n) \in \mathcal{F} \cap \mathbb{Q}[\mathbf{x}, \mathbf{y}]^n$ , the following holds:*

- i) *There exists a probabilistic algorithm that computes a solution for the weak-version of Problem (RRC) within*

$$O^{\sim} \left( \binom{t + \mathfrak{D}}{t} 8^t n^{2t+1} D^{2nt+n+2t+1} \right).$$

*arithmetic operations in  $\mathbb{Q}$  in case of success.*

- ii) *There exists a probabilistic algorithm that returns the formulas of a collection of semi-algebraic sets solving Problem (RRC) within*

$$O^{\sim} \left( \binom{t + \mathfrak{D}}{t} 8^t n^{2t+1} D^{3nt+2(n+t)+1} \right)$$

*arithmetic operations in  $\mathbb{Q}$  in case of success.*

- iii) *The semi-algebraic descriptions output by the above algorithm involves polynomials in  $\mathbb{Q}[\mathbf{y}]$  of degree bounded by  $\mathfrak{D}$ .*

We note that the binomial coefficient  $\binom{t+\mathfrak{D}}{t}$  is bounded from above by  $\mathfrak{D}^t \simeq n^t D^{nt+t}$ . Thus, the complexities given in the items i) and ii) of Theorem 5.1.3 can be bounded respectively by

$$O^{\sim}(8^t n^{3t} D^{3nt}) \text{ and } O^{\sim}(8^t n^{3t} D^{4nt}).$$

Both Theorems 5.1.2 and 5.1.3 provide a complexity  $D^{O(nt)}$  for the real root classification problem. To the best of our knowledge, there is no established complexity analysis carried out for algorithms used in the state-of-the-art software. Those algorithms rely on computing a CAD in the parameter space  $\mathbb{R}^t$  to obtain the semi-algebraic formulas. This step of computing CAD alone leads to a complexity which is doubly exponential in  $t$ .

Note that the complexity results given Theorem 5.1.3 provides a complexity with smaller constant in the exponent comparing with Theorem 5.1.2. Moreover, the degree of output polynomials in Theorem 5.1.3 is of order  $O(D^n)$  comparing to  $O(D^{2n})$  of Theorem 5.1.2.

**Organization of the chapter.** First, we present a dedicated algorithm for computing at least one point per connected component of a semi-algebraic defined by a list of inequations in Section 5.2. This algorithm and its complexity result are used in Step (b) of our algorithms. In Section 5.3, we discuss an algorithm based on Sturm-Habicht sequences for solving real root classification problem. This provides an overview on the drawbacks and potential improvements of this approach. Section 5.4 lies the definition and some useful properties of parametric Hermite matrices. There, we also present an algorithm with some optimizations to compute such a matrix. In Section 5.5, we describe our algorithm for solving the real root classification problem using this parametric Hermite matrix. The complexity analysis of the algorithms mentioned above is given in Section 5.6. Finally, in Section 5.7, we report on the practical behavior of our algorithms and illustrate its practical capabilities.

## 5.2 Computing sample points in semi-algebraic sets defined by the non-vanishing of polynomials

In this section, we study the following algorithmic problem. Given  $(g_1, \dots, g_m)$  in  $\mathbb{Q}[y_1, \dots, y_t]$ , compute at least one sample point per connected component of the semi-algebraic set  $S \subset \mathbb{R}^t$  defined by

$$g_1 \neq 0, \dots, g_m \neq 0.$$

Such sample points will be encoded with zero-dimensional parametrizations.

The main result of this section which will be used in the sequel of this paper is the following.

**Theorem 5.2.1.** *Let  $(g_1, \dots, g_m)$  in  $\mathbb{Q}[y_1, \dots, y_t]$  with  $D \geq \max_{1 \leq i \leq m} \deg(g_i)$  and  $S \subset \mathbb{R}^t$  be the semi-algebraic set defined by*

$$g_1 \neq 0, \dots, g_m \neq 0.$$

*There exists a probabilistic algorithm, which we name SamplePoints, which on input  $(g_1, \dots, g_m)$  outputs a finite family of zero-dimensional parametrizations  $\mathcal{R}_1, \dots, \mathcal{R}_k$ , all of them of degree bounded by  $(2D)^t$ , which encode at most  $(2mD)^t$  points such that  $\cup_{i=1}^k Z(\mathcal{R}_i)$  meets every connected component of  $S$  using*

$$O\left(\binom{D+t}{t} m^{t+1} 8^t D^{2t+1}\right)$$

arithmetic operations in  $\mathbb{Q}$ .

The rest of this section is devoted to the proof of this theorem.

*Proof.* By [64, Lemma 1], there exists a non-empty Zariski open set  $\mathcal{A} \times \mathcal{E} \subset \mathbb{C}^m \times \mathbb{C}$  such that for  $(\mathbf{a} = (a_1, \dots, a_m), e) \in \mathcal{A} \times \mathcal{E} \cap \mathbb{R}^m \times \mathbb{R}$ , the following holds. For  $\mathcal{I} = \{i_1, \dots, i_\ell\} \subset \{1, \dots, m\}$  and  $\sigma = (\sigma_1, \dots, \sigma_m) \in \{-1, 1\}^m$ , the algebraic sets  $V_{\mathbf{a}, e}^{\mathcal{I}, \sigma} \subset \mathbb{C}^t$  defined by

$$g_{i_1} + \sigma_{i_1} a_{i_1} e = \dots = g_{i_\ell} + \sigma_{i_\ell} a_{i_\ell} e = 0$$

are, either empty, or  $(t - \ell)$ -equidimensional and smooth, and the ideal generated by their defining equations is radical.

Note that by the transfer principle (see, e.g., Section 4.3, [9, Theorem 2.98]), one can choose instead of a scalar  $e$  an infinitesimal  $\varepsilon$  so that the algebraic sets  $V_{\mathbf{a}, \varepsilon}^{\mathcal{I}, \sigma}$  and their defining set of equations satisfy the above properties. When, in the above equations, one leaves  $\varepsilon$  as a variable, one obtains equations defining an algebraic set in  $\mathbb{C}^{t+1}$ . We denote by  $\mathfrak{V}_{\mathbf{a}, \varepsilon}^{\mathcal{I}, \sigma}$  the union of the  $(t + 1 - \ell)$ -equidimensional components of this algebraic set.

Further we also assume that the  $a_i$ 's are chosen positive.

Denote by  $\mathcal{S}^{(\varepsilon)}$  the extension of the semi-algebraic set  $\mathcal{S}$  to  $\mathbb{R}\langle\varepsilon\rangle^t$ ; similarly, the extension of any connected component  $C$  of  $\mathcal{S}$  to  $\mathbb{R}\langle\varepsilon\rangle^t$  is denoted by  $C^{(\varepsilon)}$ .

Now, remark that any connected component  $C^{(\varepsilon)}$  of  $\mathcal{S}^{(\varepsilon)}$  contains a connected component of the semi-algebraic set  $\mathcal{S}_{\mathbf{a}}^{(\varepsilon)}$  defined by:

$$(-a_1\varepsilon \geq g_1 \vee g_1 \geq a_1\varepsilon) \wedge \dots \wedge (-a_m\varepsilon \geq g_m \vee g_m \geq a_m\varepsilon)$$

Hence, we are led to compute sample points per connected component of  $\mathcal{S}_{\mathbf{a}}^{(\varepsilon)}$ . These will be encoded with zero-dimensional parametrizations with coefficients in  $\mathbb{Q}[\varepsilon]$ .

By [9, Proposition 13.1], in order to compute sample points per connected component in  $\mathcal{S}_{\mathbf{a}}^{(\varepsilon)}$ , it suffices to compute sample points in the real algebraic sets  $V_{\mathbf{a}, \varepsilon}^{\mathcal{I}, \sigma} \cap \mathbb{R}^t$ . To do that, since the algebraic sets  $V_{\mathbf{a}, \varepsilon}^{\mathcal{I}, \sigma}$  satisfy the above regularity properties, we can use the algorithm and geometric results of [171]. To state these results, one needs to introduce some notation.

Let  $\mathfrak{Q}$  be a real field,  $\mathfrak{R}$  be a real closure of  $\mathfrak{Q}$  and  $\mathfrak{C}$  be an algebraic closure of  $\mathfrak{R}$ . For an algebraic set  $V \subset \mathfrak{C}^t$  defined by  $h_1 = \dots = h_\ell = 0$  ( $h_i \in \mathfrak{Q}[\mathbf{y}]$  with  $\mathbf{y} = (y_1, \dots, y_t)$ ) and  $M \in \text{GL}(t, \mathfrak{R})$ , we denote by  $V^M$  the set  $\{M^{-1} \cdot \mathbf{x} \mid \mathbf{x} \in V\}$  and, for  $1 \leq i \leq \ell$ , by  $h_i^M$  the polynomial  $h_i(M \cdot \mathbf{y})$  and by  $\pi_i$  the canonical projection  $(y_1, \dots, y_t) \mapsto (y_1, \dots, y_i)$  ( $\pi_0$  will simply denote  $(y_1, \dots, y_t) \mapsto \{\bullet\}$ ). By slightly abusing notation, we will also denote by  $\pi_i$  projections from  $\mathfrak{V}_{\mathbf{a}, \varepsilon}^{\mathcal{I}, \sigma}$  to the first  $i$  coordinates  $(y_1, \dots, y_i)$ .

We will consider the set of critical points of the restriction of  $\pi_i$  to  $V$  and will denote this set by  $\text{crit}(\pi_i, V)$  for  $1 \leq i \leq \ell$ . Assume that  $V$  is smooth and equidimensional, by [171, Theorem 2], for a generic choice of  $M \in \text{GL}(t, \mathfrak{R})$ , the union of  $V^M \cap \pi_{t-\ell}^{-1}(0)$  with the sets  $\text{crit}(\pi_i, V^M) \cap \pi_{i-1}^{-1}(0)$  (for  $1 \leq i \leq t - \ell$ ) is finite and meets all connected components of  $V^M \cap \mathfrak{R}^t$ . Because  $V$  satisfies

the aforementioned regularity assumptions,  $\text{crit}(\pi_i, V^M) \cap \pi_{i-1}^{-1}(0)$  is defined as the projection on the  $\mathbf{y}$ -space of the solution set to the polynomials

$$\mathbf{h}^M, \quad (\lambda_1, \dots, \lambda_\ell).jac(\mathbf{h}^M, i), \quad u_1\lambda_1 + \dots + u_\ell\lambda_\ell = 1, \quad y_1 = \dots = y_{i-1} = 0,$$

where  $\mathbf{h} = (h_1, \dots, h_\ell)$ ,  $\lambda_1, \dots, \lambda_\ell$  are new variables (called Lagrange multipliers),  $jac(\mathbf{h}^M, i)$  is the Jacobian matrix associated to  $\mathbf{h}^M$  truncated by forgetting its first  $i$  columns and the  $u_i$ 's are generically chosen (see also [174, App. B]).

When  $D$  denotes the maximum degree of the  $h_j$ 's and let  $E$  be the length of a straight-line program evaluating  $\mathbf{h}$ . Observe now that, setting the  $y_j$ 's to 0 (for  $1 \leq j \leq i-1$ ), and using [175, Theorem 1] combined with the degree estimates in [175, Section 5], we obtain that such systems can be solved using

$$O\left(\left(\binom{t-i}{\ell} D^\ell (D-1)^{t-(i-1)-\ell}\right)^2 (E + (t+\ell)D + (t+\ell)^2)(t+\ell)\right)$$

arithmetic operations in  $\mathfrak{Q}$  and have at most

$$\binom{t-i}{\ell} D^\ell (D-1)^{t-(i-1)-\ell}$$

solutions.

Going back to our initial problem, one then needs to solve polynomial systems which encode the set  $\text{crit}(\pi_i, V_{\mathbf{a}, \varepsilon}^{\mathcal{I}, \sigma})$  of critical points of the restriction of  $\pi_i$  to  $V_{\mathbf{a}, \varepsilon}^{\mathcal{I}, \sigma}$ . Note that these systems have coefficients in  $\mathbb{Q}[\varepsilon]$ . To solve such systems, we rely on [178], which consists in specializing  $\varepsilon$  to a generic value  $v \in \mathbb{Q}$  and compute a zero-dimensional parametrization of the solution set to the obtained system (within the above arithmetic complexity over  $\mathbb{Q}$ ) and next use Hensel lifting and rational reconstruction to deduce from this parametrization a zero-dimensional parametrization with coefficients in  $\mathbb{Q}(\varepsilon)$ . By [178, Corollary 1] and multi-homogeneous bounds on the degree of the critical points of  $\pi_i$  to  $\mathfrak{V}_{\mathbf{a}, \varepsilon}^{\mathcal{I}, \sigma}$  as in [175, Section 5], this lifting step has a cost

$$O\left(\left((t+\ell)^4 + (t+\ell+1)E\right) \left(\binom{t-i}{\ell} D^\ell (D-1)^{t-(i-1)-\ell}\right)^2\right).$$

Hence, all in all computing one zero-dimensional parametrization for one critical locus uses

$$O\left(\left((t+\ell)^4 D + (t+\ell+1)E\right) \left(\binom{t-i}{\ell} D^\ell (D-1)^{t-(i-1)-\ell}\right)^2\right)$$

arithmetic operations in  $\mathbb{Q}$ . Note that, following [178], the degrees in  $\varepsilon$  of the numerators and denominators of the coefficients of these parametrizations are bounded by  $\binom{t}{\ell} D^\ell (D-1)^{t-\ell}$ .

Summing up for all critical loci and using

$$\sum_{i=0}^{t-\ell} \binom{t-i}{\ell} = \binom{t+1}{\ell+1},$$

the computation for a fixed  $V_{\mathbf{a},\varepsilon}^{\mathcal{I},\sigma}$  uses

$$O^{\sim} \left( ((t+\ell)^4 D + (t+\ell+1)E) \binom{t+1}{\ell+1}^2 \left( D^\ell (D-1)^{t-\ell} \right)^2 \right)$$

arithmetic operations in  $\mathbb{Q}$ . Also, the number of points computed this way is dominated by

$$\binom{t+1}{\ell+1} \left( D^\ell (D-1)^{t-\ell} \right).$$

Note that the above quantity is upper bounded by  $(2D)^t$  and bounds the degree of the output zero-dimensional parametrizations.

Taking the sum for all possible algebraic sets  $V_{\mathbf{a},\varepsilon}^{\mathcal{I},\sigma}$  and remarking that

- the sum of number of indices of cardinality  $\ell$  for  $0 \leq \ell \leq t$  is bounded by  $m^t$ ,
- the number of sets  $\sigma$  for a given  $\ell$  is bounded by  $2^t$ ,
- the sum  $\sum_{\ell=0}^t \binom{t+1}{\ell+1}^2$  equals  $2 \binom{2t+1}{t} - 1$ ,

one deduces that all these zero-dimensional parametrizations can be computed within

$$O^{\sim} \left( m^t 2^t \binom{2t+1}{t} \left( (2t)^4 D + (2t+1)E \right) D^{2t} \right)$$

arithmetic operations in  $\mathbb{Q}$  (recall that  $E$  bounds the length of a straight line program evaluating all the polynomials defining our semi-algebraic set  $\mathcal{S}$ ) which we simplify to

$$O^{\sim} (E m^t 8^t D^{2t+1}).$$

Similarly, using the above simplifications, the total number of points encoded by these zero-dimensional parametrizations is bounded above by  $(2mD)^t$ .

At this stage, we have just obtained zero-dimensional parametrizations with coefficients in  $\mathbb{Q}(\varepsilon)$ .

The above bound on the number of returned points is done but it remains to show how to specialize  $\varepsilon$  in order to get sample points per connected components in  $\mathcal{S}$ . To do that, given a parametrization  $\mathcal{R}_\varepsilon = (w, v_1, \dots, v_t) \subset \mathbb{Q}(\varepsilon)[u]^{t+1}$ , we need to find a specialization value  $e$  for  $\varepsilon$  to obtain a parametrization  $\mathcal{R}_e$  such that

- the number of real roots of the zero set associated to  $\mathcal{R}_e$  is the same as the number of real roots of the zero set associated to  $\mathcal{R}_\varepsilon$ ;
- when  $\eta$  ranges over the interval  $]0, e]$  the signs of the  $g_i$ 's at the zero set associated to  $\eta$  does not vary.

To do that, it suffices to choose  $e$  such that it is smaller than the smallest positive root of the resultant associated to  $(w, \frac{\partial w}{\partial u})$  and the smallest positive roots of the resultant associated to  $w$  and  $g_i \left( \frac{v_1}{\partial w / \partial u}, \dots, \frac{v_t}{\partial w / \partial u} \right)$ . The algebraic cost (i.e. the resultant computations) are dominated by the complexity estimates of the previous step.

Finally, note that  $E$  can be bounded by  $s \binom{D+t}{t}$  when the  $g_i$ 's are given in an expanded form in the monomial basis. Therefore, the arithmetic complexity for computing sample points of the semi-algebraic set defined by  $g_1 \neq 0, \dots, g_m \neq 0$  can be bounded by

$$O \sim \left( \binom{D+t}{t} m^{t+1} 8^t D^{2t+1} \right).$$

□

**Remark 5.2.2.** Observe that the coefficients of the rational parametrizations with coefficients in  $\mathbb{Q}[\varepsilon]$  have bit size depending both on the maximum bit size  $\tau$  of the coefficients of the input polynomials  $g_1, \dots, g_m$  and the bit size of the generically chosen  $a_i$ 's.

When substituting the infinitesimal  $\varepsilon$  by a small enough rational number  $e \in \mathbb{Q}$ , one obtains zero-dimensional parametrizations with coefficients in  $\mathbb{Q}$  of bit size depending on the one of  $e$  also. Admissible values for  $e$  depend on the magnitude of the real roots of the univariate resultant we exhibit in the above proof. Because we start with rational parametrizations of degree bounded by  $O(D)^t$ , assuming that the bit size of the  $a_i$ 's is bounded by  $O(D)^t$  (following the reasoning like the one in [57]), one could show using standard quantitative results that the bit size of  $e$  may be  $\tau D^{O(t)}$  (because  $e$  is obtained through the isolation of real roots of a univariate polynomial of degree  $D^{O(t)}$ ). However, this is a worst case analysis and most of the time, we observe in practice that one can choose for  $e$  values of reasonable bit size.

We end this section with a corollary which is a consequence of the proof of [9, Theorem 13.18]. Basically, once we have the parametrizations computed by the algorithm on which Theorem 5.2.1 relies, one can compute sample points per connected components of the semi-algebraic set  $\mathcal{S}$  within the same arithmetic complexity bounds. The idea is just to evaluate the  $g_i$ 's at these rational parametrizations and use bounds on the minimal distance between two roots of a univariate polynomial such as [9, Prop. 10.22]. Hence, the proof of the corollary below follows *mutatis mutandis* the same steps as the one of [9, Theorem 13.18].

**Corollary 5.2.3.** Let  $(g_1, \dots, g_m)$  in  $\mathbb{Q}[y_1, \dots, y_t]$  with  $D \geq \max_{1 \leq i \leq m} \deg(g_i)$  and  $\mathcal{S} \subset \mathbb{R}^t$  be the semi-algebraic set defined by

$$g_1 \neq 0, \dots, g_m \neq 0.$$

There exists a probabilistic algorithm which on input  $(g_1, \dots, g_m)$  outputs a finite set of points  $\mathcal{P}$  in  $\mathbb{Q}^t$  of cardinality at most  $(2mD)^t$  points such that  $\mathcal{P}$  meets every connected component of  $\mathcal{S}$  using

$$O \sim \left( \binom{D+t}{t} m^{t+1} 8^t D^{2t+1} \right).$$

arithmetic operations in  $\mathbb{Q}$ .

Note that, by contrast with Theorem 5.2.1, the above corollary shows how to obtain output points with coordinates in  $\mathbb{Q}$ .

### 5.3 Algorithm based on Sturm-Habicht sequences

In this section, we describe an algorithm based on Sturm-Habicht sequences for solving Problem (RRC) and discuss its shortcomings.

We consider a sequence  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  where  $\mathbf{y} = (y_1, \dots, y_t)$  and  $\mathbf{x} = (x_1, \dots, x_n)$ . Let  $D$  be an upper bound of the total degree of the  $f_i$ 's. We require that the input system  $\mathbf{f}$  satisfies the properties below.

**Assumption 5.B.** *Let  $\mathbf{f}$  be the above parametric polynomial system and  $\mathcal{V}$  be the algebraic set defined by  $\mathbf{f}$ . We say that  $\mathbf{f}$  satisfies Assumptions (5.B) if the following properties hold.*

- (B1) *The ideal generated by  $\mathbf{f}$  is radical.*
- (B2) *The algebraic set  $\mathcal{V}$  is equidimensional of dimension  $t$ .*
- (B3) *The restriction of  $\pi : (\mathbf{y}, \mathbf{x}) \mapsto \mathbf{y}$  to  $\mathcal{V}$  is dominant.*

It is well-known that the above assumptions are satisfied by sufficiently generic systems (see e.g. [175]).

In what follows, we rely on the existence of a parametric geometric resolution due to Schost [178] to reduce our initial multivariate problem to a univariate one.

Using [178, Proposition 2] with Assumption (B1), there exists a non-empty open Zariski set  $\mathcal{A}$  of  $\mathbb{C}^n$  such that, for  $(a_1, \dots, a_n) \in \mathbb{Q}^n \cap \mathcal{A}$ , there exists a parametric geometric resolution  $(w_{\mathbf{a}}, v_1, \dots, v_n) \in \mathbb{Q}(\mathbf{y})[u]^n$  of  $\mathbf{f}$  that satisfies the following properties.

- $w_{\mathbf{a}}$  is a square-free polynomial in  $\mathbb{Q}[\mathbf{y}][u]$ .
- $u = \sum_{i=1}^n a_i x_i$ .
- There exists a non-empty Zariski open subset  $\mathcal{Y}_{\mathbf{a}} \subset \mathbb{C}^t$  such that, for  $\eta \in \mathcal{Y}_{\mathbf{a}}$ , we have that

$$V(\mathbf{f}(\eta, \cdot)) = \left\{ \left( \frac{v_1}{\partial w_{\mathbf{a}} / \partial u}(\eta, \vartheta), \dots, \frac{v_n}{\partial w_{\mathbf{a}} / \partial u}(\eta, \vartheta) \right) \mid w_{\mathbf{a}}(\eta, \vartheta) = 0, \frac{\partial w_{\mathbf{a}}}{\partial u}(\eta, \vartheta) \neq 0 \right\}.$$

The set  $\mathcal{Y}_{\mathbf{a}}$  can be chosen as the set where the leading coefficient of  $w_{\mathbf{a}}$ , the resultant of  $w_{\mathbf{a}}$  and  $\partial w_{\mathbf{a}} / \partial u$ , and the denominators appearing in the coefficients of  $v_1, \dots, v_n$  do not vanish.

As a consequence, for  $\eta \in \mathcal{Y}_a$ , the number of complex solutions of  $\mathbf{f}(\eta, \cdot)$  is invariant and equals the partial degree of  $w_a$  in  $u$ . We denote by  $\Delta$  the partial degree of  $w_a$  in  $u$ . By Bézout's inequality (see e.g. [101]),  $\Delta$  is bounded above by  $D^n$ .

Let  $\eta \in \mathbb{C}^t$  and  $w_a(\eta, \cdot)$  be the specialization of the  $\mathbf{y}$  variables in  $w_a$  at  $\eta$ . From the existence of such a parametric resolution, we deduce that, for  $\eta \in \mathcal{Y}_a$ , the map

$$\varphi : (x_1, \dots, x_n) \mapsto \sum_{i=1}^n a_i x_i$$

is a bijection between the complex roots of  $\mathbf{f}(\eta, \cdot)$  and  $w_a(\eta, \cdot)$ .

**Lemma 5.3.1.** *Let  $\mathbf{f}$  be a parametric system satisfying Assumption (5.B) and  $w_a$  be the eliminating polynomial in the parametric geometric resolution of  $\mathbf{f}$  as above. Then, we have*

$$V(\langle f_1, \dots, f_s, u - \sum_{i=1}^n a_i x_i \rangle \cap \mathbb{Q}[\mathbf{y}][u]) = V(w_a).$$

Consequently, the total degree of  $w_a$  is at most  $D^n$ .

*Proof.* We prove that, under Assumption (5.B), there exists a square-free polynomial  $w \in \mathbb{Q}[\mathbf{y}][x]$  satisfying

$$V(\langle f_1, \dots, f_s, u - \sum_{i=1}^n a_i x_i \rangle \cap \mathbb{Q}[\mathbf{y}][u]) = V(w).$$

Let  $\pi_u : \mathbb{C}^{t+n+1} \rightarrow \mathbb{C}^{t+1}$ ,  $(\mathbf{y}, \mathbf{x}, u) \mapsto (\mathbf{y}, u)$  and  $\mathcal{V}_u$  be the algebraic set defined by  $\langle \mathbf{f}, u - \sum_{i=1}^n a_i x_i \rangle$ . Note that  $\mathcal{V}$  and  $\mathcal{V}_u$  are isomorphic taking the map  $(\mathbf{y}, \mathbf{x}) \mapsto (\mathbf{y}, \mathbf{x}, \sum_{i=1}^n a_i x_i)$  as an isomorphism between them. Then, as the algebraic set  $\mathcal{V}$  satisfies Assumption (5.B),  $\mathcal{V}_u$  is equidimensional of dimension  $t$  and the restriction of  $\Pi : \mathbb{C}^{t+n+1} \rightarrow \mathbb{C}^t$ ,  $(\mathbf{y}, \mathbf{x}, u) \mapsto \mathbf{y}$  to  $\mathcal{V}_u$  is dominant. Therefore, the Zariski closure of  $\pi_u(\mathcal{V}_u)$  is an equidimensional algebraic set of dimension  $t$ . Hence, there exists a square-free polynomial  $w \in \mathbb{Q}[\mathbf{y}][u]$  such that  $V(w) = \overline{\pi_u(\mathcal{V}_u)}$ . Therefore, we obtain  $V(\langle f_1, \dots, f_s, u - \sum_{i=1}^n a_i x_i \rangle \cap \mathbb{Q}[\mathbf{y}][u]) = \overline{\pi_u(\mathcal{V}_u)} = V(w)$ .

It remains to show that  $w_a$  equals to  $w$  up to a constant. By the definition of parametric geometric resolution, for  $\eta \in \mathcal{Y}_a$ , then  $w_a(\eta, \cdot)$  and  $w(\eta, \cdot)$  share the same complex roots. Therefore,  $w_a$  equals to  $w$  up to a factor in  $\mathbb{Q}[\mathbf{y}]$ . However, both  $w_a$  and  $w$  do not contain such kind of factor.

By Bézout's inequalities, the degree of  $V(f_1, \dots, f_s, u - \sum_{i=1}^n a_i x_i)$  is at most  $D^n$ . Hence, the degree of  $\overline{\pi_u(\mathcal{V}_u)}$  is also bounded by  $D^n$ . Therefore, the total degree of  $w_a$  is bounded by  $D^n$ .  $\square$

Recall that  $\mathcal{Y}_a$  is the non-empty Zariski open subset of  $\mathbb{C}^t$  where the leading coefficient of  $w_a$ , the resultant of  $w_a$  and  $\partial w_a / \partial u$ , and the denominators appearing in the coefficients of  $v_1, \dots, v_n$  do not vanish. Lemma 5.3.2 shows that the numbers of real roots of  $\mathbf{f}(\eta, \cdot)$  and  $w_a(\eta, \cdot)$  also coincide over  $\mathcal{Y}_a$ .

**Lemma 5.3.2.** *Let  $\mathcal{Y}_\alpha$  be as above. Then, for  $\eta \in \mathcal{Y}_\alpha \cap \mathbb{R}^t$ , the numbers of real solutions of  $w_\alpha(\eta, \cdot)$  and  $\mathbf{f}(\eta, \cdot)$  are equal.*

*Proof.* Let  $\eta \in \mathbb{R}^t \cap \mathcal{Y}_\alpha$ . By definition of  $\mathcal{Y}_\alpha$ , the restriction of  $\varphi : (x_1, \dots, x_n) \mapsto \sum_{i=1}^n a_i x_i$  to  $V(\mathbf{f}(\eta, \cdot))$  is a bijection between the complex roots of  $\mathbf{f}(\eta, \cdot)$  and  $w_\alpha(\eta, \cdot)$ .

As the sequence  $\mathbf{f}(\eta, \cdot)$  contains polynomials of coefficients in  $\mathbb{R}$ , the non-real complex roots of  $\mathbf{f}(\eta, \cdot)$  appears as pairs of conjugate complex points of  $\mathbb{C}^n$ . Assume that there exists a non-real root whose image by  $\varphi$  is a real root of  $w_\alpha(\eta, \cdot)$ , then its conjugate is also mapped to the same real root. This contradicts the bijectivity of  $\varphi$ . Therefore, the numbers of real solutions of  $\mathbf{f}(\eta, \cdot)$  and  $w_\alpha(\eta, \cdot)$  coincide.  $\square$

For  $h \in \mathbb{Q}[\mathbf{y}][u]$  of degree  $\Delta$  in  $u$ , we denote by

$$\Sigma \left( h, \frac{\partial h}{\partial u} \right) = (s_0, \dots, s_\Delta) \subset \mathbb{Q}[\mathbf{y}]$$

the leading coefficients of the signed subresultant sequence associated to  $(h, \partial h / \partial u)$  (see [9, Notation 4.21]). Here we enumerate this sequence in a way such that  $s_0$  is the leading coefficient of  $h$  as a polynomial in  $u$  and  $s_\Delta$  is the resultant of  $h$  and  $\partial h / \partial u$ .

We recall the specialization property of signed subresultant coefficients (Proposition 4.4.12). For  $\eta \in \mathbb{R}^t$  that does not cancel the leading coefficient of  $h$  as a polynomial in  $u$ , then the signed subresultant coefficients of  $h(\eta, \cdot)$  and  $\partial h(\eta, \cdot) / \partial u$  are exactly the evaluation of  $(s_0, \dots, s_\Delta)$  at  $\eta$ .

By Theorem 4.4.17, the number of real roots of  $h(\eta, \cdot)$  equals the generalized permanences minus variations (see Definition 4.4.16) of  $(s_0, \dots, s_\Delta)_\eta$ . Note that this value is uniquely defined upon a sign pattern of  $(s_0, \dots, s_\Delta)_\eta$ .

We can now describe Algorithm 5.1 which takes as input a parametric polynomial sequence  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  satisfying Assumption (5.B) and it outputs semi-algebraic formulas solving Problem (RRC).

It uses the following subroutines:

- **EliminatingPolynomial** which takes as input  $\mathbf{f}$  and outputs an eliminating polynomial  $w_\alpha$ , i.e., the first polynomial in a parametric geometric resolution of  $\mathbf{f}$ .

Such an algorithm can be derived from the probabilistic algorithm given in [178] that computes parametric geometric resolutions.

- **SignedSubresultantCoefficients** which takes as input  $w_\alpha$  and outputs  $\Sigma(w_\alpha, \partial w_\alpha / \partial u) = (s_0, \dots, s_\Delta)$ .

We refer to [9, Algo. 8.77] for the explicit description of such an algorithm.

- **SamplePoints** which takes as input the signed subresultant coefficients  $(s_0, \dots, s_\Delta) \subset \mathbb{Q}[\mathbf{y}]$  and outputs at least one point per connected components of the semi-algebraic set defined by  $\{s_i \neq 0 \mid 1 \leq i \leq \Delta, s_i \text{ is not a constant}\}$ .

We refer to Theorem 5.2.1 in Section 5.2 for the explicit description of such an algorithm and its complexity.

- PermanencesMinusVariations which takes as input a sequence  $(s_0, \dots, s_\Delta)_\eta$  and return its generalized permanences minus variations.

The generalized permanences minus variations can be computed using its definition given in Definition 4.4.16 (see also [9, Algo. 9.4]).

---

**Algorithm 5.1:** RRC-Sturm-Habicht

---

**Input:** A parametric system  $\mathbf{f} \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  satisfying Assumption (5.B)  
**Output:** Semi-algebraic descriptions solving Problem (RRC) for the input  $\mathbf{f}$

- 1  $w_\alpha \leftarrow \text{EliminatingPolynomial}(\mathbf{f})$
- 2  $(s_0, \dots, s_\Delta) \leftarrow \text{SignedSubresultantCoefficients}(w_\alpha, \partial w_\alpha / \partial u, u)$
- 3  $L \leftarrow \text{SamplePoints}(\{s_i \neq 0 \mid 1 \leq i \leq \Delta, s_i \text{ is not a constant}\})$
- 4 **for**  $\eta \in L$  **do**
- 5      $r_\eta \leftarrow \text{PermanencesMinusVariations}((s_0, \dots, s_\Delta)_\eta)$
- 6 **end**
- 7 **return**  $\{(\text{sign}(s_0, \dots, s_\Delta)_\eta, \eta, r_\eta) \mid \eta \in L\}$

---

**Theorem 5.1.2.** *Let  $\mathbf{f} = (f_1, \dots, f_n) \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  be a parametric system and  $D$  be the largest total degree among the  $\deg(f_i)$ 's. We assume that  $\mathbf{f}$  satisfies Assumption (5.B).*

*Then, Algorithm 5.1, which is probabilistic, computes semi-algebraic formulas solving Problem RRC within*

$$O\left(\binom{t + 2D^{2n}}{t} 32^t D^{5nt+3n}\right)$$

*arithmetic operations in  $\mathbb{Q}$ . These semi-algebraic formulas contain polynomials in  $\mathbb{Q}[\mathbf{y}]$  of degree bounded by  $2D^{2n}$ .*

*Proof.* We start with the correctness statement. Recall that  $s_0$  is the leading coefficient of  $w_\alpha$  as a polynomial in  $u$ . By [9, Proposition 8.74], for  $\eta \in \mathbb{C}^t$  that does not cancel  $s_0$ , the signed subresultant coefficients of  $w_\alpha(\eta, \cdot)$  and  $\partial w_\alpha(\eta, \cdot) / \partial u$  is the specialization of  $(s_0, \dots, s_\Delta)$  at  $\eta$ . Therefore, from [9, Theorem 4.33], the number of real roots of  $w_\alpha(\eta, \cdot)$  can be derived from the sign of the sequence  $(s_0, \dots, s_\Delta)_\eta$  for  $\eta \notin V(s_0)$ .

On the other hand, the semi-algebraic set  $\mathcal{S}$  defined by

$$\{s_i \neq 0 \mid 1 \leq i \leq \Delta, s_i \text{ is not a constant}\}$$

is composed of open semi-algebraic connected components, namely  $\mathcal{S}_1, \dots, \mathcal{S}_\ell$ . Over each of them, the signed subresultant coefficients  $s_i$  are sign-invariant. Thus, the number of distinct real roots of  $w_\alpha(\eta, \cdot)$  is invariant when  $\eta$  varies in  $\mathcal{S}_i$  for each  $1 \leq i \leq \ell$ .

Recall that  $\mathcal{Y}_a \subset \mathbb{C}^t$  is the non-empty Zariski open set in Lemma 5.3.2 such that for  $\eta \in \mathcal{Y}_a$ , the numbers of real roots of  $\mathbf{f}(\eta, \cdot)$  and  $w_a(\eta, \cdot)$  coincide. Therefore, the number of real solutions of  $\mathbf{f}(\eta, \cdot)$  is also invariant when  $\eta$  varies in  $\mathcal{S}_i \cap \mathcal{Y}_a$ .

Let  $L$  be the set of sample points of  $\mathcal{S}$ . We deduce from the above arguments that the semi-algebraic sets defined by

$$\bigwedge_{i=1}^{\Delta} \text{sign}(s_i) = \text{sign}(s_i(\eta))$$

for  $\eta \in L$  satisfy the requirement of Problem RRC. The correctness of our algorithm is proven.

By Lemma 5.3.1, the total degree of  $w_a$  is bounded above by  $d^n$ . Using [9, Proposition 8.71] on the polynomial  $w_a$  and  $\partial w_a / \partial u$ , we obtain the bound

$$\deg s_j \leq D^n(2j - 1) \leq 2D^n$$

for  $0 \leq j \leq \Delta$ . Using this bound, we are now able to analyze the complexity of Algorithm 5.1.

By [178, Corollary 1], running EliminatingPolynomial on input  $\mathbf{f} = (f_1, \dots, f_n)$  where the total degree of each  $f_i$  is bounded by  $d$  takes

$$O^{\sim} \left( \binom{4D^n + t}{t} D^n \right)$$

arithmetic operations in  $\mathbb{Q}$ .

The signed subresultant coefficients of  $w_a$  and  $\partial w_a / \partial u$  can be computed using an evaluation-interpolation scheme as follows.

As the degree of  $s_i$  is bounded by  $2D^{2n}$ , we need to compute the signed subresultant coefficients of the evaluation of  $(w_a, \partial w_a / \partial u)$  at  $\binom{t+2D^{2n}}{t}$  distinct points. Note that  $\binom{t+2D^{2n}}{t}$  is bounded by  $2^t D^{2nt}$ .

Using [9, Algo. 8.77], it yields an arithmetic complexity  $O(D^{2n})$  for each of those signed subresultant computations. Hence, in total, the specialized signed subresultant coefficients can be computed by

$$O(2^t D^{2nt+2n})$$

arithmetic operations in  $\mathbb{Q}$ .

Next, the cost of interpolating the  $s_i$ 's can be bounded by

$$O^{\sim}(\Delta 2^t D^{2nt}) \simeq O^{\sim}(2^t D^{2nt+n})$$

using the interpolation given in [36]. Thus, the arithmetic complexity of SignedSubresultantCoefficients lies in

$$O^{\sim}(2^t D^{2nt+n}).$$

We rely on Corollary 5.2.3 for estimating the complexity of SamplePoints. Using the algorithm of Section 5.2 (see Theorem 5.2.1 and Corollary 5.2.3) on the sequence  $(s_0, \dots, s_{\Delta})$ ,

one can compute sample points per connected components of the semi-algebraic set defined by  $\{s_i \neq 0 \mid 1 \leq i \leq \Delta, s_i \text{ is not a constant}\}$  in time

$$O\left(\binom{t + 2D^{2n}}{t} t^4 D^{nt+n} 8^t (2D^{2n})^{2t+1}\right) \simeq O\left(\binom{t + 2D^{2n}}{t} 32^t D^{5nt+3n}\right).$$

By Corollary 5.2.3, this subroutine outputs a finite subset of  $\mathbb{Q}^t$  whose cardinal is bounded by  $4^t D^{3nt}$ . Using [9, Algorithm 9.4] to compute the permanences minus variations leads to an arithmetic complexity of

$$O(4^t D^{3nt+n}).$$

Summing up all the partial costs, we conclude that Algorithm 5.1 runs within

$$O\left(\binom{t + 2D^{2n}}{t} 32^t D^{5nt+3n}\right)$$

arithmetic operations in  $\mathbb{Q}$ . □

**Example 5.3.3.** *We will illustrate the algorithms of this paper using the sequence*

$$\mathbf{f} = (x_1^2 + x_2^2 - y_1, x_1x_2 + y_2x_2 + y_3x_1).$$

*We choose  $u = x_2$  when running Algorithm 5.1 (in a reasonable implementation, one would pick randomly a linear form but we choose this one to obtain smaller data).*

*We obtain the following rational parametrization  $(w, v_1, v_2)$  with*

$$\begin{aligned} w &= u^4 + 2y_3u^3 + (y_2^2 + y_3^2 - y_1)u^2 - 2y_1y_3u - y_1y_3^2, \\ v_2 &= 2y_3u^3 + (2y_2^2 + 2y_3^2 - 2y_1)u^2 - 6y_1y_3u - 4y_1y_3^2, \\ v_1 &= 2y_2u^3 + 2y_1y_3y_2. \end{aligned}$$

*The signed subresultant coefficients associated to  $(w, \frac{\partial w}{\partial u})$  are:*

$$\begin{aligned} s_0 &= 1, s_1 = 1, s_2 = -2y_2^2 + y_3^2 + 2y_1, \\ s_3 &= -y_2^6 - 2y_2^4y_3^2 - y_2^2y_3^4 + 3y_1y_2^4 - 14y_1y_2^2y_3^2 + y_1y_3^4 - 3y_1^2y_2^2 - 2y_1^2y_3^2 + y_1^3, \\ s_4 &= (y_2y_3)^2y_1(-y_2^6 - 3y_2^4y_3^2 - 3y_2^2y_3^4 - y_3^6 + 3y_1y_2^4 - 21y_1y_2^2y_3^2 + 3y_1y_3^4 - 3y_1^2y_2^2 - 3y_1^2y_3^2 + y_1^3). \end{aligned}$$

*Since  $s_0$  and  $s_1$  are constants, we then compute at least one point per connected component of the semi-algebraic set defined by*

$$s_2 \neq 0 \wedge s_3 \neq 0 \wedge s_4 \neq 0.$$

*This is done using e.g. RAGlib (the Real Algebraic Geometry library) [170] (because it implements an algorithm which is easy to use for sampling points in semi-algebraic sets). We obtain 35 points and find that the realizable sign conditions for  $(s_2, s_3, s_4)$  are*

$$[-1, -1, -1], [-1, -1, 1], [-1, 1, 1], [1, -1, -1], [1, -1, 1], [1, 1, -1], [1, 1, 1].$$

Applying [9, Theorem 4.32], we deduce the corresponding numbers of real roots to these sign patterns

$$\begin{aligned} 0 \text{ real root} &\rightarrow (s_2 < 0 \wedge s_3 < 0 \wedge s_4 > 0) \vee (s_2 < 0 \wedge s_3 > 0 \wedge s_4 > 0) \vee (s_2 > 0 \wedge s_3 < 0 \wedge s_4 > 0), \\ 2 \text{ real roots} &\rightarrow (s_2 < 0 \wedge s_3 < 0 \wedge s_4 < 0) \vee (s_2 > 0 \wedge s_3 < 0 \wedge s_4 < 0) \vee (s_2 > 0 \wedge s_3 > 0 \wedge s_4 < 0), \\ 4 \text{ real roots} &\rightarrow (s_2 > 0 \wedge s_3 > 0 \wedge s_4 > 0). \end{aligned}$$

Note that the maximum degree of the polynomials involved in the above formulas is 11. By contrast, observe that the restriction of the projection  $\pi : (\mathbf{x}, \mathbf{y}) \rightarrow \mathbf{y}$  to the real algebraic set defined by  $\mathbf{f}$  is proper. Hence, applying a semi-algebraic variant of Thom's isotopy lemma as in [21], one can deduce that the set of critical values of this map discriminates the regions of the parameters' space over which the number of real roots of  $\mathbf{f}$  remains invariant.

Using immediate Gröbner bases computations, one obtains that the Zariski closure of this set of critical values is defined by the vanishing of

$$y_1(-y_2^6 - 3y_2^4y_3^2 - 3y_2^2y_3^4 - y_3^6 + 3y_1y_2^4 - 21y_1y_2^2y_3^2 + 3y_1y_3^4 - 3y_1^2y_2^2 - 3y_1^2y_3^2 + y_1^3)$$

which has only degree 7.

## 5.4 Parametric Hermite matrices

In this section, we adapt the construction of Hermite matrices to the context of parametric systems and describe an algorithm for computing those *parametric Hermite matrices*. We will use these matrices to obtain a better algorithm than the one presented in Section 5.3 for real root classification.

### 5.4.1 Definition

In this subsection, we present the definition of our parametric Hermite matrix and prove some properties which will be used further in this chapter.

Let  $\mathbf{f} = (f_1, \dots, f_s)$  be a polynomial sequence in  $\mathbb{Q}[\mathbf{y}][\mathbf{x}]$ . We take the rational function field  $\mathbb{Q}(\mathbf{y})$  as the base field  $\mathbb{K}$  and denote by  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  the ideal generated by  $\mathbf{f}$  in  $\mathbb{K}[\mathbf{x}]$ . We require that the system  $\mathbf{f}$  satisfies Assumption (5.A).

This leads to the following lemma, which is the foundation for the construction of our parametric Hermite matrices.

**Lemma 5.4.1.** *Assume that  $\mathbf{f}$  satisfies Assumption (5.A). Then the ideal  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  is zero-dimensional.*

*Proof.* Assume that there exists a coordinate  $x_i$  for  $1 \leq i \leq n$  such that  $\langle \mathbf{f} \rangle \cap \mathbb{C}[\mathbf{y}, x_i] = \langle 0 \rangle$ . We denote respectively by  $\pi_i$  and  $\tilde{\pi}_i$  the projections  $(\mathbf{y}, \mathbf{x}) \mapsto (\mathbf{y}, x_i)$  and  $(\mathbf{y}, x_i) \mapsto \mathbf{y}$ . By the assumption above,  $\overline{\pi_i(\mathcal{V})}$  is the whole space  $\mathbb{C}^{t+1}$ . Then, we have the identity

$$\mathbb{C}^{t+1} = \overline{(\tilde{\pi}_i^{-1}(\mathcal{O}) \cup \tilde{\pi}_i^{-1}(\mathbb{C}^t \setminus \mathcal{O})) \cap \pi_i(\mathcal{V})},$$

where  $\mathcal{O}$  be the dense Zariski open subset of  $\mathbb{C}^t$  required in Assumption (5.A).

Since  $\tilde{\pi}_i$  is a map from  $\mathbb{C}^{t+1}$  to  $\mathbb{C}^t$ , its fibers are of dimension at most 1. Therefore, we have that  $\dim \tilde{\pi}_i^{-1}(\mathbb{C}^t \setminus \mathcal{O}) \leq 1 + \dim(\mathbb{C}^t \setminus \mathcal{O}) \leq t$ . As Assumption (5.A) holds and  $\dim \tilde{\pi}_i^{-1}(\mathbb{C}^t \setminus \mathcal{O}) \leq t$ , we have that  $\dim \tilde{\pi}_i^{-1}(\mathcal{O}) \cap \pi_i(\mathcal{V}) = t$ . This contradicts to the above identity above. We conclude that, for  $1 \leq i \leq n$ ,  $\langle \mathbf{f} \rangle \cap \mathbb{C}[\mathbf{y}, x_i] \neq \langle 0 \rangle$ .

On the other hand, by Assumption (5.A), the Zariski-closure of  $\pi(\mathcal{V})$  is the whole parameter space  $\mathbb{C}^t$ . Thus, we have that  $\langle \mathbf{f} \rangle \cap \mathbb{C}[\mathbf{y}] = \langle 0 \rangle$ . Since  $\langle \mathbf{f} \rangle \cap \mathbb{C}[\mathbf{y}] = (\langle \mathbf{f} \rangle \cap \mathbb{C}[\mathbf{y}, x_i]) \cap \mathbb{C}[\mathbf{y}]$  for every  $1 \leq i \leq n$ , there exists a polynomial  $p_i \in \langle \mathbf{f} \rangle \cap \mathbb{C}[\mathbf{y}, x_i]$  whose degree with respect to  $x_i$  is non-zero. Clearly,  $p_i$  is an element of the ideal  $\langle \mathbf{f} \rangle_{\mathbb{K}}$ . Thus, there exists  $d_i$  such that  $x_i^{d_i}$  is a leading term in  $\langle \mathbf{f} \rangle_{\mathbb{K}}$ . Hence,  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  is a zero-dimensional ideal.  $\square$

Lemma 5.4.1 allows us to apply the construction of Hermite matrices described in Section 4.4.3 to parametric systems as follows.

Since the ideal  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  is zero-dimensional by Lemma 5.4.1, its associated quotient ring  $A_{\mathbb{K}} = \mathbb{K}[\mathbf{x}] / \langle \mathbf{f} \rangle_{\mathbb{K}}$  is a finite dimensional  $\mathbb{K}$ -vector space (Theorem 3.3.1). Let  $\delta$  denote the dimension of  $A_{\mathbb{K}}$  as a  $\mathbb{K}$ -vector space.

We consider a basis  $B = \{b_1, \dots, b_\delta\}$  of  $A_{\mathbb{K}}$ , where the  $b_i$ 's are taken as monomials in the variables  $\mathbf{x}$ . Such a basis can be derived from Gröbner bases as follows. We fix an admissible monomial ordering  $\succ$  over the set of monomials in the variables  $\mathbf{x}$  and compute a Gröbner basis  $G$  with respect to the ordering  $\succ$  of the ideal  $\langle \mathbf{f} \rangle_{\mathbb{K}}$ . Then, the monomials that are not divisible by any leading monomial of elements of  $G$  form a basis of  $A_{\mathbb{K}}$ .

For an element  $p \in \mathbb{K}[\mathbf{x}]$ , we denote by  $\bar{p}$  the class of  $p$  in the quotient ring  $A_{\mathbb{K}}$ . A representative of  $\bar{p}$  can be derived by computing the normal form of  $p$  by the Gröbner basis  $G$ , which results in a linear combination of elements of  $B$  with coefficients in  $\mathbb{Q}(\mathbf{y})$ . The map of multiplication by  $p$ , denoted by  $\mathcal{L}_p$ , is an endomorphism of  $A_{\mathbb{K}}$  defined as

$$\begin{aligned} \mathcal{L}_p : A_{\mathbb{K}} &\rightarrow A_{\mathbb{K}}, \\ \bar{q} &\mapsto \overline{p \cdot q}. \end{aligned}$$

Recall that Hermite's quadratic form of the ideal  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  is defined as the bilinear form that sends  $(\bar{p}, \bar{q}) \in A_{\mathbb{K}} \times A_{\mathbb{K}}$  to the trace of  $\mathcal{L}_{p \cdot q}$ .

Assume now the basis  $B$  of  $A_{\mathbb{K}}$  is fixed. Every multiplication map  $\mathcal{L}_p$  admits a matrix representation with respect to  $B$ , whose entries are elements in  $\mathbb{Q}(\mathbf{y})$ . The trace of  $\mathcal{L}_p$  can be computed as the trace of the matrix representing it. Similarly, Hermite's quadratic form of  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  can be represented by a matrix with respect to  $B$ . This leads to the following definition.

**Definition 5.4.2.** *Given a parametric polynomial system  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  satisfying Assumption (5.A). Let  $\mathbb{K} = \mathbb{Q}(\mathbf{y})$ , we fix a basis  $B = \{b_1, \dots, b_\delta\}$  of the vector space  $\mathbb{K}[\mathbf{x}] / \langle \mathbf{f} \rangle_{\mathbb{K}}$ . The parametric Hermite matrix of  $\mathbf{f}$  with respect to the basis  $B$  is defined as the symmetric matrix*

$$\mathcal{H} = (h_{i,j})_{1 \leq i,j \leq \delta},$$

where  $h_{i,j} = \text{trace}(\mathcal{L}_{b_i \cdot b_j})$ .

It is important to note that the definition of parametric Hermite matrices depends both on the input system  $\mathbf{f}$  and the choice of the monomial basis  $B$ .

**Example 5.4.3.** We consider the same system  $\mathbf{f} = (x_1^2 + x_2^2 - y_1, x_1x_2 + y_2x_2 + y_3x_1)$  as in Example 5.3.3. The parametric Hermite matrix  $\mathcal{H}_1$  associated to  $\mathbf{f}$  with respect to the basis  $B_1 = \{1, x_2, x_1, x_2^2\}$  is

$$\begin{bmatrix} 4 & -2y_3 & -2y_2 & -2(y_2^2 + y_3^2 + y_1) \\ * & -2(y_2^2 + y_3^2 + y_1) & 4y_2y_3 & 2(3y_2^2y_3 - y_3^3) \\ * & * & 2(y_2^2 - y_3^2 + y_1) & 2(y_2^3 - 3y_2y_3^2 - y_1y_2) \\ * & * & * & 2y_2^4 - 12y_2^2y_3^2 + 2y_3^4 - 4y_1y_2^2 + 2y_1^2 \end{bmatrix}.$$

Whereas, using the lexicographical ordering  $x_1 \succ x_2$ , we obtain the basis  $B_2 = \{1, x_2, x_2^2, x_2^3\}$ . The matrix  $\mathcal{H}_2$  associated to  $\mathbf{f}$  with respect to  $B_2$  is the following Hankel matrix:

$$\begin{bmatrix} 4 & -2y_3 & -2y_2^2 + 2y_3^2 + 2y_1 & 6y_2^2y_3 - 2y_3^3 \\ * & * & * & 2y_2^4 - 12y_2^2y_3^2 + 2y_3^4 - 4y_1y_2^2 + 2y_1^2 \\ * & * & * & -10y_2^4y_3 + 20y_2^2y_3^3 - 2y_3^5 + 10y_1y_2^2y_3 \\ * & * & * & -2y_2^6 + 30y_2^4y_3^2 - 30y_2^2y_3^4 + 2y_3^6 + 6y_1y_2^4 - 18y_1y_2^2y_3^2 - 6y_1^2y_2^2 + 2y_1^3 \end{bmatrix}.$$

We remark that all the entries of the matrices above lie in  $\mathbb{Q}[\mathbf{y}]$  and that the entries of the second matrix are of higher degree than the first one's. Moreover, writing the basis change from  $B_1$  to  $B_2$ , we obtain the factorization of  $\mathcal{H}_2 = P^T \cdot \mathcal{H}_1 \cdot P$  where

$$P = \begin{bmatrix} 1 & 0 & 0 & y_1y_3 \\ 0 & 1 & 0 & -y_2^2 + y_1 \\ 0 & 0 & 0 & -y_2y_3 \\ 0 & 0 & 1 & -y_3 \end{bmatrix}.$$

## 5.4.2 Gröbner bases and parametric Hermite matrices

In the previous subsection, we have defined parametric Hermite matrices assuming one knows a Gröbner basis  $G$  with respect to some monomial ordering of the ideal  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  where  $\mathbb{K} = \mathbb{Q}(\mathbf{y})$  and  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  is the ideal of  $\mathbb{K}[\mathbf{x}]$  generated by  $\mathbf{f}$ .

Computing such a Gröbner basis may be costly as this would require to perform arithmetic operations over the field  $\mathbb{Q}(\mathbf{y})$  (or  $\mathbb{Z}/p\mathbb{Z}(\mathbf{y})$  where  $p$  is a prime when tackling this computational task through modular computations). In this paragraph, we show that one can obtain parametric Hermite matrices by considering some Gröbner bases of the ideal  $\langle \mathbf{f} \rangle \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  (hence, enabling the use of efficient implementations of Gröbner basis algorithms such as the F4/F5 algorithms [59, 60]).

Since the graded reverse lexicographical ordering (*grevlex* for short) is known for yielding Gröbner bases of relatively small degree comparing to other orders, we prefer using this ordering to construct our parametric Hermite matrices. Further, we will use the notation  $\text{grevlex}(\mathbf{x})$  for the *grevlex* ordering among the variables  $\mathbf{x}$  (with  $x_1 \succ \dots \succ x_n$ ) and  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$  (with  $y_1 \succ \dots \succ y_t$ ) for the elimination ordering. We denote respectively by  $\text{lm}_{\mathbf{x}}(p)$  and  $\text{lc}_{\mathbf{x}}(p)$  the leading monomial and the leading coefficient of  $p \in \mathbb{K}[\mathbf{x}]$  with respect to the ordering  $\text{grevlex}(\mathbf{x})$ .

**Lemma 5.4.4.** *Let  $\mathcal{G}$  be the reduced Gröbner basis of  $\langle \mathbf{f} \rangle$  with respect to the elimination ordering  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$ . Then  $\mathcal{G}$  is also a Gröbner basis of  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  with respect to the ordering  $\text{grevlex}(\mathbf{x})$ .*

*Proof.* Since  $\mathcal{G}$  is a Gröbner basis of the ideal  $\langle \mathbf{f} \rangle$ , every polynomial  $f_i$  of  $\mathbf{f}$  can be written as  $f_i = \sum_{g \in \mathcal{G}} c_g \cdot g$  where  $c_g \in \mathbb{Q}[\mathbf{x}, \mathbf{y}]$ . Therefore, any element of  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  can also be written as a combination of elements of  $\mathcal{G}$  with coefficients in  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$ . In other words,  $\mathcal{G}$  is a set of generators of  $\langle \mathbf{f} \rangle_{\mathbb{K}}$ .

Let  $p$  be a polynomial in  $\mathbb{K}[\mathbf{x}]$ ,  $p$  is contained in  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  if and only if there exists a polynomial  $q \in \mathbb{Q}[\mathbf{y}]$  such that  $q \cdot p \in \langle \mathbf{f} \rangle$ . Thus, the leading monomial of  $p$  as an element of  $\mathbb{K}[\mathbf{x}]$  with respect to the grevlex ordering  $\text{grevlex}(\mathbf{x})$  is contained in the ideal  $\langle \text{lm}_{\mathbf{x}}(g) \mid g \in \mathcal{G} \rangle$ . Therefore,  $\mathcal{G}$  is a Gröbner basis of  $\langle \mathbf{f} \rangle_{\mathbb{K}}$ .  $\square$

Hereafter, we denote by  $\mathcal{G}$  the reduced Gröbner basis of  $\langle \mathbf{f} \rangle$  with respect to the elimination ordering  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$ . Let  $\mathcal{B}$  be the set of all monomials in  $\mathbf{x}$  that are not reducible by  $\mathcal{G}$ , which is finite by Lemmas 5.4.1 and 5.4.4. The set  $\mathcal{B}$  actually forms a basis of the  $\mathbb{K}$ -vector space  $\mathbb{K}[\mathbf{x}]/\langle \mathbf{f} \rangle_{\mathbb{K}}$ . Then, we denote by  $\mathcal{H}$  the parametric Hermite matrix associated to  $\mathbf{f}$  with respect to this basis  $\mathcal{B}$ .

We consider the following assumption on the input system  $\mathbf{f}$ .

**Assumption 5.C.** *For  $g \in \mathcal{G}$ , the leading coefficient  $\text{lc}_{\mathbf{x}}(g)$  does not depend on the parameters  $\mathbf{y}$ .*

As the computations in the quotient ring  $A_{\mathbb{K}}$  are done through normal form reductions by  $\mathcal{G}$ , the lemma below is straight-forward.

**Lemma 5.4.5.** *Under Assumption (5.C), the entries of the parametric Hermite matrix  $\mathcal{H}$  are elements of  $\mathbb{Q}[\mathbf{y}]$ .*

*Proof.* Since Assumption (5.C) holds, the leading coefficients  $\text{lc}_{\mathbf{x}}(g)$  do not depend on parameters  $\mathbf{y}$  for  $g \in \mathcal{G}$ . The normal form reduction in  $A_{\mathbb{K}}$  of any polynomial in  $\mathbb{Q}[\mathbf{y}][\mathbf{x}]$  returns a polynomial in  $\mathbb{Q}[\mathbf{y}][\mathbf{x}]$ . Thus, each normal form can be written as a linear combination of  $\mathcal{B}$  whose coefficients lie in  $\mathbb{Q}[\mathbf{y}]$ . Hence, the multiplication map  $\mathcal{L}_{b_i \cdot b_j}$  for  $1 \leq i, j \leq \delta$  can be represented by polynomial matrices in  $\mathbb{Q}[\mathbf{y}]$  with respect to the basis  $\mathcal{B}$ . As an immediate consequence, the entries of  $\mathcal{H}$ , as being the traces of those multiplication maps, are polynomials in  $\mathbb{Q}[\mathbf{y}]$ .  $\square$

The next proposition states that Assumption (5.C) is satisfied by a generic system  $\mathbf{f}$ . It implies that the entries of the parametric Hermite matrix of a generic system with respect to the basis  $\mathcal{B}$  derived from  $\mathcal{G}$  completely lie in  $\mathbb{Q}[\mathbf{y}]$ . We postpone the proof of Proposition 5.4.6 to Subsection 5.6.1 where we prove a more general result (see Proposition 5.6.1).

**Proposition 5.4.6.** *Let  $\mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}$  be the set of polynomials in  $\mathbb{C}[\mathbf{x}, \mathbf{y}]$  having total degree bounded by  $D$ . There exists a non-empty Zariski open subset  $\mathcal{F}_C$  of  $\mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^n$  such that Assumption (5.C) is satisfied by any  $\mathbf{f} \in \mathcal{F}_C \cap \mathbb{Q}[\mathbf{x}, \mathbf{y}]^n$  with  $\mathcal{G}$  the Gröbner basis of  $\mathbf{f}$  w.r.t. the ordering  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$ .*

### 5.4.3 Specialization property of parametric Hermite matrices

Recall that  $\mathcal{G}$  is the reduced Gröbner basis of  $\langle \mathbf{f} \rangle$  with respect to the ordering  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$  and  $\mathcal{B}$  is the basis of  $\mathbb{K}[\mathbf{x}]/\langle \mathbf{f} \rangle_{\mathbb{K}}$  derived from  $\mathcal{G}$  as discussed in the previous subsection. Then,  $\mathcal{H}$  is the parametric Hermite matrix associated to  $\mathbf{f}$  with respect to the basis  $\mathcal{B}$ .

Let  $\eta \in \mathbb{C}^t$  and  $\phi_\eta : \mathbb{C}(\mathbf{y})[\mathbf{x}] \rightarrow \mathbb{C}[\mathbf{x}]$ ,  $p(\mathbf{y}, \mathbf{x}) \mapsto p(\eta, \mathbf{x})$  be the specialization map that evaluates the parameters  $\mathbf{y}$  at  $\eta$ . Then  $\mathbf{f}(\eta, \cdot) = (\phi_\eta(f_1), \dots, \phi_\eta(f_s))$ . We denote by  $\mathcal{H}(\eta)$  the specialization  $(\phi_\eta(h_{i,j}))_{1 \leq i, j \leq \delta}$  of  $\mathcal{H}$  at  $\eta$ .

Recall that, for a polynomial  $p \in \mathbb{C}(\mathbf{y})[\mathbf{x}]$ , the leading coefficient of  $p$  considered as a polynomial in the variables  $\mathbf{x}$  with respect to the ordering  $\text{grevlex}(\mathbf{x})$  is denoted by  $\text{lc}_{\mathbf{x}}(p)$ . In this subsection, for  $p \in \mathbb{C}[\mathbf{x}]$ , we use  $\text{lm}(p)$  to denote the leading monomial of  $p$  with respect to the ordering  $\text{grevlex}(\mathbf{x})$ .

Let  $\mathcal{W}_\infty \subset \mathbb{C}^t$  denote the algebraic set  $\cup_{g \in \mathcal{G}} V(\text{lc}_{\mathbf{x}}(g))$ . In Proposition 5.4.8, we prove that, outside  $\mathcal{W}_\infty$ , the specialization  $\mathcal{H}(\eta)$  coincides with the classical Hermite matrix of the zero-dimensional ideal  $\mathbf{f}(\eta, \cdot) \subset \mathbb{Q}[\mathbf{x}]$ . This is the main result of this subsection.

Since the operations over the  $\mathbb{K}$ -vector space  $A_{\mathbb{K}}$  rely on normal form reductions by the Gröbner basis  $\mathcal{G}$  of  $\langle \mathbf{f} \rangle_{\mathbb{K}}$ , the specialization property of  $\mathcal{H}$  depends on the specialization property of  $\mathcal{G}$ . Lemma 5.4.7 below, which is a direct consequence of [118, Theorem 3.1], provides the specialization property of  $\mathcal{G}$ . We give here a more elementary proof for this lemma than the one in [118].

**Lemma 5.4.7.** *Let  $\eta \in \mathbb{C}^t \setminus \mathcal{W}_\infty$ . Then the specialization  $\mathcal{G}(\eta, \cdot) := \{\phi_\eta(g) \mid g \in \mathcal{G}\}$  is a Gröbner basis of the ideal  $\langle \mathbf{f}(\eta, \cdot) \rangle \subset \mathbb{C}[\mathbf{x}]$  generated by  $\mathbf{f}(\eta, \cdot)$  with respect to the ordering  $\text{grevlex}(\mathbf{x})$ .*

*Proof.* Since  $\eta \in \mathbb{C}^t \setminus \mathcal{W}_\infty$ , the leading coefficient  $\text{lc}_{\mathbf{x}}(g)$  does not vanish at  $\eta$  for every  $g \in \mathcal{G}$ . Thus,  $\text{lm}_{\mathbf{x}}(g) = \text{lm}(\phi_\eta(g))$ .

We denote by  $\mathcal{M}$  the set of all monomials in the variables  $\mathbf{x}$  and

$$\mathcal{M}_{\mathcal{G}} := \{m \in \mathcal{M} \mid \exists g \in \mathcal{G} : \text{lm}_{\mathbf{x}}(g) \mid m\} = \{m \in \mathcal{M} \mid \exists g \in \mathcal{G} : \text{lm}(\phi_\eta(g)) \mid m\}.$$

For any  $p \in \langle \mathbf{f} \rangle \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$ , we prove that  $\text{lm}(\phi_\eta(p)) \in \mathcal{M}_{\mathcal{G}}$ . If  $p$  is identically zero, there is nothing to prove. So, we assume that  $p \neq 0$ ,  $p$  is then expanded in the form below:

$$p = \sum_{m \in \mathcal{M}_{\mathcal{G}}} c_m \cdot m + \sum_{m \in \mathcal{M} \setminus \mathcal{M}_{\mathcal{G}}} c_m \cdot m,$$

where the  $c_m$ 's are elements of  $\mathbb{Q}[\mathbf{y}]$ . Since  $p$  is not identically zero, there exists  $m \in \mathcal{M}_{\mathcal{G}}$  such that  $c_m \neq 0$ .

Since  $\mathcal{G}$  is a Gröbner basis of  $\langle \mathbf{f} \rangle_{\mathbb{K}}$ , any monomial in  $\mathcal{M}_{\mathcal{G}}$  can be reduced by  $\mathcal{G}$  to a unique normal form in  $\mathbb{K}[\mathbf{x}]$ . These divisions involve denominators, which are products of some powers of the leading coefficients of  $\mathcal{G}$  with respect to the variables  $\mathbf{x}$ . We write

$$\text{NF}_{\mathcal{G}}(p) = \sum_{m \in \mathcal{M}_{\mathcal{G}}} c_m \cdot \text{NF}_{\mathcal{G}}(m) + \sum_{m \in \mathcal{M} \setminus \mathcal{M}_{\mathcal{G}}} c_m \cdot m.$$

As  $p \in \langle \mathbf{f} \rangle_{\mathbb{K}}$ , we have that  $\text{NF}_{\mathcal{G}}(p) = 0$ , which implies

$$\sum_{m \in \mathcal{M} \setminus \mathcal{M}_{\mathcal{G}}} c_m \cdot m = - \sum_{m \in \mathcal{M}_{\mathcal{G}}} c_m \cdot \text{NF}_{\mathcal{G}}(m).$$

Therefore, we have the identity

$$p = \sum_{m \in \mathcal{M}_{\mathcal{G}}} c_m \cdot (m - \text{NF}_{\mathcal{G}}(m))$$

Since  $\eta$  does not cancel any denominator appearing in  $\text{NF}_{\mathcal{G}}(m)$ , we can specialize the identity above without any problem:

$$\phi_{\eta}(p) = \sum_{m \in \mathcal{M}_{\mathcal{G}}} \phi_{\eta}(c_m) \cdot (m - \phi_{\eta}(\text{NF}_{\mathcal{G}}(m))).$$

If at least one of the  $\phi_{\eta}(c_m)$  does not vanish, then the leading monomial of  $\phi_{\eta}(f)$  is in  $\mathcal{M}_{\mathcal{G}}$ . Otherwise, if all the  $\phi_{\eta}(c_m)$  are canceled, then  $\phi_{\eta}(p)$  is identically zero, and there is not any new leading monomial appearing either. So, the leading monomial of any  $p \in \langle \mathbf{f}_{\eta} \rangle$  is contained in  $\mathcal{M}_{\mathcal{G}}$ , which means  $\mathcal{G}(\eta, \cdot)$  is a Gröbner basis of  $\langle \mathbf{f}(\eta, \cdot) \rangle$  with respect to  $\text{grevlex}(\mathbf{x})$ .  $\square$

**Proposition 5.4.8.** *For any  $\eta \in \mathbb{C}^t \setminus \mathcal{W}_{\infty}$ , the specialization  $\mathcal{H}(\eta)$  coincides with the classic Hermite matrix of the zero-dimensional ideal  $\langle \mathbf{f}(\eta, \cdot) \rangle \subset \mathbb{C}[\mathbf{x}]$ .*

*Proof.* As a consequence of Lemma 5.4.7, each computation in  $A_{\mathbb{K}}$  derives a corresponding one in  $\mathbb{C}[\mathbf{x}]/\langle \mathbf{f}(\eta, \cdot) \rangle$  by evaluating  $\mathbf{y}$  at  $\eta$  in every normal form reduction by  $\mathcal{G}$ . This evaluation is allowed since  $\eta$  does not cancel any denominator appearing during the computation. Therefore, we deduce immediately the specialization property of the Hermite matrix.  $\square$

Using Proposition 5.4.8 and [9, Theorem 4.102], we obtain immediately the following corollary that allows us to use parametric Hermite matrices to count the root of a specialization of a parametric system.

**Corollary 5.4.9.** *Let  $\eta \in \mathbb{C}^t \setminus \mathcal{W}_{\infty}$ , then the rank of  $H(\eta)$  is the number of distinct complex roots of  $\mathbf{f}(\eta, \cdot)$ . When  $\eta \in \mathbb{R}^t \setminus \mathcal{W}_{\infty}$ , the signature of  $H(\eta)$  is the number of distinct real roots of  $\mathbf{f}(\eta, \cdot)$ .*

*Proof.* By Proposition 5.4.8,  $\mathcal{H}(\eta)$  is a Hermite matrix of the zero-dimensional ideal  $\langle \mathbf{f}(\eta, \cdot) \rangle$ . Then, [9, Theorem 4.102] implies that the rank (resp. the signature) of  $\mathcal{H}(\eta)$  equals to the number of distinct complex (resp. real) solutions of  $\mathbf{f}(\eta, \cdot)$ .  $\square$

We finish this subsection by giving some explanation for what happens above  $\mathcal{W}_{\infty}$ , where our parametric Hermite matrix  $\mathcal{H}$  does not have good specialization property.

**Lemma 5.4.10.** *Let  $\mathcal{W}_{\infty}$  be defined as above. Then  $\mathcal{W}_{\infty}$  contains the following sets:*

- The non-proper points of the restriction of  $\pi$  to  $\mathcal{V}$ .
- The set of points  $\eta \in \mathbb{C}^t$  such that the fiber  $\pi^{-1}(\eta) \cap \mathcal{V}$  is infinite.
- The image by  $\pi$  of the irreducible components of  $\mathcal{V}$  whose dimensions are smaller than  $t$ .

*Proof.* The claim for the set of non-properness of the restriction of  $\pi$  to  $\mathcal{V}$  is already proven in [134, Theorem 2]. We focus on the two remaining sets.

Using the Hermite matrix, we know that for  $\eta \in \mathbb{C}^t \setminus \mathcal{W}_\infty$ , the system  $\mathbf{f}(\eta, \cdot)$  admits a non-empty finite set of complex solutions. On the other hand, for any  $\eta \in \mathbb{C}^t$  such that  $\pi^{-1}(\eta) \cap \mathcal{V}$  is infinite,  $\mathbf{f}(\eta, \cdot)$  has infinitely many complex solutions. Therefore, the set of such points  $\eta$  is contained in  $\mathcal{W}_\infty$ .

Let  $\mathcal{V}_{>t}$  be the union of irreducible components of  $\mathcal{V}$  of dimension greater than  $t$ . By the fiber dimension theorem [184, Theorem 1.25], the fibers of the restriction of  $\pi$  to  $\mathcal{V}_{>t}$  must have dimension at least one. Similarly, the components of dimension  $t$  whose images by  $\pi$  are contained in a Zariski closed subset of  $\mathbb{C}^t$  also yield infinite fibers. Therefore, as proven above, all of these components are contained in  $\pi^{-1}(\mathcal{W}_\infty)$ .

We now consider the irreducible components of dimension smaller than  $t$ . Let  $\mathcal{V}_{\geq t}$  and  $\mathcal{V}_{<t}$  be respectively the union of irreducible components of  $\mathcal{V}$  of dimension at least  $t$  and at most  $t-1$ . We have that  $\mathcal{V} = \mathcal{V}_{\geq t} \cup \mathcal{V}_{<t}$ . Let  $I \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  denote the ideal generated by  $\mathbf{f}$ . Using the primary decomposition of  $I$  (see e.g. [48, Sec. 4.8]), we have that  $I$  is the intersection of two ideals  $I_{\geq t}$  and  $I_{<t}$  such that  $V(I_{\geq t}) = \mathcal{V}_{\geq t}$  and  $V(I_{<t}) = \mathcal{V}_{<t}$ . We write

$$I = I_{\geq t} \cap I_{<t}.$$

We denote by  $R$  the polynomial ring  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$ . Then, the above identity is transferred into  $R$ :

$$I \cdot R = (I_{\geq t} \cdot R) \cap (I_{<t} \cdot R).$$

Since  $\dim(\overline{\pi(\mathcal{V}_{<t})}) \leq t-1$ , then there exists a non-zero polynomial  $p \in I_{<t} \cap \mathbb{Q}[\mathbf{y}]$ . As  $p$  is a unit in  $\mathbb{Q}(\mathbf{y})$ , the ideal  $I_{<t} \cdot R$  is exactly  $R$ . So,

$$I \cdot R = I_{\geq t} \cdot R.$$

Note that, by Lemma 5.4.4,  $\mathcal{G}$  is a Gröbner basis of  $I \cdot R$ , then it is also a Gröbner basis of  $I_{\geq t} \cdot R$ . Therefore, the Hermite matrices associated to  $I$  and  $I_{\geq t}$  (with respect to the basis derived from  $\mathcal{G}$ ) coincide. So, for  $\eta \notin \mathcal{W}_\infty$ , the ranks of those matrices are equal and so are the numbers of complex points in  $\pi^{-1}(\eta) \cap \mathcal{V}$  and  $\pi^{-1}(\eta) \cap \mathcal{V}_{\geq t}$ . As  $\pi^{-1}(\eta) \cap \mathcal{V}_{\geq t} \subset \pi^{-1}(\eta) \cap \mathcal{V}$ , we have that  $\pi^{-1}(\eta) \cap \mathcal{V} = \pi^{-1}(\eta) \cap \mathcal{V}_{\geq t}$ . This leads to

$$\pi^{-1}(\mathbb{C}^t \setminus \mathcal{W}_\infty) \cap \mathcal{V}_{\geq t} = \pi^{-1}(\mathbb{C}^t \setminus \mathcal{W}_\infty) \cap \mathcal{V}.$$

Then,  $\pi^{-1}(\mathbb{C}^t \setminus \mathcal{W}_\infty) \cap \mathcal{V}_{<t} = \emptyset$  or equivalently,  $\mathcal{V}_{<t} \subset \pi^{-1}(\mathcal{W}_\infty)$ , which concludes the proof.  $\square$

#### 5.4.4 Computing parametric Hermite matrices

Given  $\mathbf{f} = (f_1, \dots, f_s) \in \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  satisfying Assumption (5.A). We keep on denoting  $\mathbb{K} = \mathbb{Q}(\mathbf{y})$ . Let  $\mathcal{G}$  be the reduced Gröbner basis of  $\langle \mathbf{f} \rangle$  with respect to the ordering  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$  and  $\mathcal{B}$  be the set of all monomials in the variables  $\mathbf{x}$  which are not reducible by  $\mathcal{G}$ . The set  $\mathcal{B}$  then forms a basis of the  $\mathbb{K}$ -vector space  $\mathbb{K}[\mathbf{x}]/\langle \mathbf{f} \rangle_{\mathbb{K}}$ .

In this subsection, we focus on the computation of the parametric Hermite matrix associated to  $\mathbf{f}$  with respect to the basis  $\mathcal{B}$ .

Note that one can design an algorithm using only the definition of parametric Hermite matrices given in Subsection 5.4.1. More precisely, for each  $b_i \cdot b_j \in \mathcal{B}$  ( $1 \leq i, j \leq \delta$ ), one computes the matrix representing  $\mathcal{L}_{b_i \cdot b_j}$  in the basis  $\mathcal{B}$  by computing the normal form of every  $b_i \cdot b_j \cdot b_k$  for  $1 \leq k \leq \delta$ . Therefore, in total, this direct algorithm requires  $O(\delta^3)$  normal form reductions of polynomials in  $\mathbb{K}[\mathbf{x}]$ .

In Algorithm 5.2 below, we present another algorithm for computing  $\mathcal{H}$ . It uses the following subroutines:

- **GrobnerBasis** that takes as input the system  $\mathbf{f}$  and computes the reduced Gröbner basis  $\mathcal{G}$  of  $\langle \mathbf{f} \rangle$  with respect to the ordering  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$  and the basis  $\mathcal{B} = \{b_1, \dots, b_\delta\} \subset \mathbb{Q}[\mathbf{x}]$  of  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]/\langle \mathbf{f} \rangle_{\mathbb{Q}(\mathbf{y})}$  derived from  $\mathcal{G}$ .

Such an algorithm can be obtained using any general algorithm for computing Gröbner basis, which we refer to F4/F5 algorithms [59, 60].

- **ReduceGB** that takes as input the Gröbner basis  $\mathcal{G}$  and outputs a subset  $\mathcal{G}'$  of  $\mathcal{G}$  which is still a Gröbner basis of  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  with respect to the ordering  $\text{grevlex}(\mathbf{x})$ .

This subroutine aims to remove the elements in  $\mathcal{G}$  that we do not need. Even though  $\mathcal{G}$  is reduced as a Gröbner basis of  $\langle \mathbf{f} \rangle$  with respect to  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$ , it is not necessarily the reduced Gröbner basis of  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  with respect to  $\text{grevlex}(\mathbf{x})$ . Using [48, Lemma 3, Sec. 2.7], we can design ReduceGB to remove all the elements of  $\mathcal{G}$  which have duplicate leading monomials (in  $\mathbf{x}$ ). We obtain as output a subset  $\mathcal{G}'$  of  $\mathcal{G}$  which is also a Gröbner basis  $\mathcal{G}'$  for  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  with respect to  $\text{grevlex}(\mathbf{x})$ . Note that this tweak reduces not only the cardinality of the Gröbner basis in use but also the size of the set  $\mathcal{W}_\infty$  introduced in Subsection 5.4.3 (as we have less leading coefficients).

- **XMatrices** that takes as input  $(\mathcal{G}', \mathcal{B})$  and computes the matrix representation of the multiplication maps  $\mathcal{L}_{x_i}$  ( $1 \leq i \leq n$ ) with respect to  $\mathcal{B}$ .

This computation is done directly by reducing every  $x_i \cdot b_j$  ( $1 \leq i \leq n, 1 \leq j \leq \delta$ ) to its normal form in  $\mathbb{K}[\mathbf{x}]/\langle \mathbf{f} \rangle_{\mathbb{K}}$  using  $\mathcal{G}'$ .

- **BMatrices** that takes as input the matrices representing  $(\mathcal{L}_{x_1}, \dots, \mathcal{L}_{x_n})$  and  $\mathcal{B}$  and computes the matrices representing the  $\mathcal{L}_{b_i}$ 's ( $1 \leq i \leq \delta$ ) in the basis  $\mathcal{B}$ .

We design BMatrices in a way that it constructs the matrices of  $\mathcal{L}_{b_i}$ 's inductively in the degree of the  $b_i$ 's as follows.

At the beginning, we have the multiplication matrices of 1 and the  $x_i$ 's; those are the matrices of the elements of degree zero and one. Note that, for any element  $b$  of  $\mathcal{B}$ . At the step of computing the matrix of an element  $b \in \mathcal{B}$ , we remark that there exist a variable  $x_i$  and a monomial  $b' \in \mathcal{B}$  such that  $b = x_i \cdot b'$  and the matrix of  $b'$  is already computed (as  $\deg(b') < \deg(b)$ ). Therefore, we simply multiply the matrices of  $\mathcal{L}_{x_i}$  and  $\mathcal{L}_{b'}$  to obtain the matrix of  $\mathcal{L}_b$ .

- **TraceComputing** that takes as input the multiplication matrices  $\mathcal{L}_{b_1}, \dots, \mathcal{L}_{b_\delta}$  and computes the matrix  $(\text{trace}(\mathcal{L}_{b_i \cdot b_j}))_{1 \leq i \leq j \leq \delta}$ . This matrix is actually the parametric Hermite matrix  $\mathcal{H}$  associated to  $\mathbf{f}$  with respect to the basis  $\mathcal{B}$ . To design this subroutine, we use the following remark given in [166].

Let  $p, q \in \mathbb{K}[\mathbf{x}]$ . The normal form  $\bar{p}$  of  $p$  by  $\mathcal{G}$  can be written as  $\bar{p} = \sum_{i=1}^{\delta} c_i \cdot b_i$  where the  $c_i$ 's lie in  $\mathbb{K}$ . Then, we have the identity

$$\text{trace}(\mathcal{L}_{p \cdot q}) = \sum_{i=1}^{\delta} c_i \cdot \text{trace}(\mathcal{L}_{q \cdot b_i}),$$

Hence, by choosing  $p = b_i \cdot b_j$  and  $q = 1$ , we can compute  $h_{i,j}$  using the normal form  $\overline{b_i \cdot b_j}$  and  $\text{trace}(\mathcal{L}_{b_1}), \dots, \text{trace}(\mathcal{L}_{b_\delta})$ .

Note that  $\text{trace}(\mathcal{L}_{b_i})$  is easily computed from the matrix of the map  $\mathcal{L}_{b_i}$ . On the other hand, the normal form  $\overline{b_i \cdot b_j}$  can be read off from the  $j$ -th row of the matrix representing  $\mathcal{L}_{b_i}$ , which is already computed at this point.

It is also important to notice that there are many duplicated entries in  $\mathcal{H}$ . Thus, we should avoid all the unnecessary re-computation. This is done easily by keeping a list tracking distinct entries of  $\mathcal{H}$ .

The pseudo-code of Algorithm 5.2 is presented below. Its correctness follows simply from our definition of parametric Hermite matrices.

Besides the parametric Hermite matrix  $\mathcal{H}$ , we return a polynomial  $\mathbf{w}_\infty$  which is the square-free part of  $\text{lcm}_{g \in \mathcal{G}}(\text{lc}_x(g))$  for further usage. Note that  $V(\mathbf{w}_\infty) = \mathcal{W}_\infty$ .

---

**Algorithm 5.2:** DRL-HermiteMatrix

---

**Input:** A parametric polynomial system  $\mathbf{f} = (f_1, \dots, f_s)$

**Output:** A parametric Hermite matrix  $\mathcal{H}$  associated to  $\mathbf{f}$  with respect to the basis  $\mathcal{B}$

- 1  $\mathcal{G}, \mathcal{B} \leftarrow \text{GröbnerBasis}(\mathbf{f}, \text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y}))$
  - 2  $\mathcal{G}' \leftarrow \text{ReduceGB}(\mathcal{G})$
  - 3  $\mathbf{w}_\infty \leftarrow \text{sqfree}(\text{lcm}_{g \in \mathcal{G}}(\text{lc}_x(g)))$
  - 4  $(\mathcal{L}_{x_1}, \dots, \mathcal{L}_{x_n}) \leftarrow \text{XMatrices}(\mathcal{G}', \mathcal{B})$
  - 5  $(\mathcal{L}_{b_1}, \dots, \mathcal{L}_{b_\delta}) \leftarrow \text{BMatrices}((\mathcal{L}_{x_1}, \dots, \mathcal{L}_{x_n}), \mathcal{B})$
  - 6  $\mathcal{H} \leftarrow \text{TraceComputing}(\mathcal{L}_{b_1}, \dots, \mathcal{L}_{b_\delta})$
  - 7 **return**  $[\mathcal{H}, \mathbf{w}_\infty]$
-

**Removing denominators.** Note that, through the computation in the quotient ring  $A_{\mathbb{K}}$ , the entries of our parametric Hermite matrix possibly contains denominators that lie in  $\mathbb{Q}[\mathbf{y}]$ . As the algorithm that we introduce in Section 5.5 will require us to manipulate the parametric Hermite matrix that we compute, these denominators can be a bottleneck to handle the matrix. Therefore, we introduce an extra subroutine `RemoveDenominator` that returns a parametric Hermite matrix  $\mathcal{H}'$  of  $\mathbf{f}$  without denominator.

- `RemoveDenominator` that takes as input the matrix  $\mathcal{H}$  computed by `DRL-HermiteMatrix` and outputs a matrix  $\mathcal{H}'$  which is the parametric Hermite matrix associated to  $\mathbf{f}$  with respect to a basis  $\mathcal{B}'$  that will be made explicit below.

As we can freely choose any basis of the form  $\{c_i \cdot b_i \mid 1 \leq i \leq \delta\}$  where the  $c_i$ 's are elements of  $\mathbb{Q}[\mathbf{y}]$ , we should use a basis that leads to a denominator-free matrix. To do this, we choose  $c_i$  as the denominator of  $\text{trace}(\mathcal{L}_{b_i})$  (which lies in the first row of the matrix  $\mathcal{H}$  computed by `TraceComputing`). Then, for the entry of  $\mathcal{H}$  that corresponds to  $b_i$  and  $b_j$ , we can multiply it with  $c_i \cdot c_j$ . The output matrix  $\mathcal{H}'$  is the parametric Hermite matrix associated to  $\mathbf{f}$  with respect to the basis  $\{c_i \cdot b_i \mid 1 \leq i \leq \delta\}$ .

We observe in many examples that this subroutine returns either a denominator-free matrix or a matrix with smaller degree denominators. Thus, it facilitates further computations on the output matrix.

**Evaluation & interpolation scheme for generic systems.** Here we assume that the input system  $\mathbf{f}$  satisfies Assumption (5.C). By Lemma 5.4.5, the entries of  $\mathcal{H}$  are polynomials in  $\mathbb{Q}[\mathbf{y}]$ . Suppose that we know beforehand a value  $\Lambda$  that is larger than the degree of any entry of  $\mathcal{H}$ , we can compute  $\mathcal{H}$  by an evaluation & interpolation scheme as follows.

We start by choosing randomly a set  $\mathcal{E}$  of  $\binom{t+\Lambda}{t}$  distinct points in  $\mathbb{Q}^t$ . Then, for each  $\eta \in \mathcal{E}$ , we use `DRL-HermiteMatrix` (Algorithm 5.2) on the input  $\mathbf{f}(\eta, \cdot)$  to compute the classical Hermite matrix associated to  $\mathbf{f}(\eta, \cdot)$  with respect to the ordering  $\text{grevlex}(\mathbf{x})$ . These computations involve only polynomials in  $\mathbb{Q}[\mathbf{x}]$  and not in  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$ . Finally, we interpolate the parametric Hermite matrix  $\mathcal{H}$  from its specialized images  $\mathcal{H}(\eta)$  computed previously.

Since Assumption (5.C) holds, then  $\mathcal{W}_\infty$  is empty. By Proposition 5.4.8, the Hermite matrix of  $\mathbf{f}(\eta, \cdot)$  with respect to  $\text{grevlex}(\mathbf{x})$  is the image  $\mathcal{H}(\eta)$  of  $\mathcal{H}$ . Therefore, the above scheme computes correctly the parametric Hermite matrix  $\mathcal{H}$ .

We also remark that, in the computation of the specializations  $\mathcal{H}(\eta)$ , we can replace the subroutine `XMatrices` in `DRL-HermiteMatrix` by a linear-algebra-based algorithm described in [58]. That algorithm constructs the Macaulay matrix and carries out matrix reductions to obtain simultaneously the normal forms that `XMatrices` requires.

In Section 5.6, we will estimate the complexity of this evaluation & interpolation scheme when the input system  $\mathbf{f}$  satisfies some generic assumptions.

## 5.5 Algorithms for real root classification

We present in this section two algorithms targeting the real root classification problem through parametric Hermite matrices. The one described in Subsection 5.5.1 aims to solve the weak version of Problem (RRC). The second algorithm, given in Subsection 5.5.2 outputs the semi-algebraic formulas of the cells  $\mathcal{S}_i$  that solves Problem (RRC). Further, in Section 5.6, we will see that, for a generic sequence  $\mathbf{f}$ , the semi-algebraic formulas computed by this algorithm consist of polynomials of degree bounded by  $n(D-1)D^n$ , which is better than the degree bound  $2D^{2n}$  obtained by Algorithm 5.1 and all previously known bounds.

Throughout this section, our input is a parametric polynomial system  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$ . We require that  $\mathbf{f}$  satisfies Assumptions (5.A) and that the ideal  $\langle \mathbf{f} \rangle$  is radical.

Let  $\mathcal{G}$  be the reduced Gröbner basis of the ideal  $\langle \mathbf{f} \rangle \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  with respect to the ordering  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$ . Let  $\mathbb{K}$  denote the rational function field  $\mathbb{Q}(\mathbf{y})$ . We recall that  $\mathcal{B} \subset \mathbb{Q}[\mathbf{x}]$  is the basis of  $\mathbb{K}[\mathbf{x}]/\langle \mathbf{f} \rangle_{\mathbb{K}}$  derived from  $\mathcal{G}$  and  $\mathcal{H}$  is the parametric Hermite matrix associated to  $\mathbf{f}$  with respect to the basis  $\mathcal{B}$ .

### 5.5.1 Algorithm for weak real root classification

From Subsection 5.4.3, we know that, outside the algebraic set  $\mathcal{W}_{\infty} := \cup_{g \in \mathcal{G}} V(\text{lc}_{\mathbf{x}}(g))$ , the parametric matrix  $\mathcal{H}$  possesses good specialization properties (see Proposition 5.4.8). We denote by  $\mathbf{w}_{\infty}$  the square-free part of  $\text{lcm}_{g \in \mathcal{G}} \text{lc}_{\mathbf{x}}(g)$ . This polynomial  $\mathbf{w}_{\infty}$  is returned as an output of Algorithm 5.2. Note that  $V(\mathbf{w}_{\infty}) = \mathcal{W}_{\infty}$ .

**Lemma 5.5.1.** *When Assumption (5.A) holds and the ideal  $\langle \mathbf{f} \rangle$  is radical, the determinant of  $\mathcal{H}$  is not identically zero.*

*Proof.* Recall that  $\mathbb{K}$  denotes the rational function field  $\mathbb{Q}(\mathbf{y})$ . We prove that the ideal  $\langle \mathbf{f} \rangle_{\mathbb{K}} \subset \mathbb{K}[\mathbf{x}]$  is radical.

Let  $p \in \mathbb{K}[\mathbf{x}]$  such that there exists  $k \in \mathbb{N}$  satisfying  $p^k \in \langle \mathbf{f} \rangle_{\mathbb{K}}$ . Therefore, there exists a polynomial  $q \in \mathbb{Q}[\mathbf{y}]$  such that  $q \cdot p^k \in \langle \mathbf{f} \rangle$ . Then,  $(q \cdot p)^k \in \langle \mathbf{f} \rangle$ . As  $\langle \mathbf{f} \rangle$  is radical, we have that  $q \cdot p \in \langle \mathbf{f} \rangle$ . Thus,  $p \in \langle \mathbf{f} \rangle_{\mathbb{K}}$ , which concludes that  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  is radical.

By Lemma 5.4.1,  $\langle \mathbf{f} \rangle_{\mathbb{K}}$  is a radical zero-dimensional ideal in  $\mathbb{Q}(\mathbf{y})$ . Since  $\mathcal{H}$  is also a Hermite matrix (in the classic sense) of  $\langle \mathbf{f} \rangle_{\mathbb{K}}$ ,  $\mathcal{H}$  is full rank. Therefore,  $\det(\mathcal{H})$  is not identically zero.  $\square$

Let  $\mathbf{w}_{\mathcal{H}} := \mathbf{n} / \text{gcd}(\mathbf{n}, \mathbf{w}_{\infty})$  where  $\mathbf{n}$  is the square-free part of the numerator of  $\det(\mathcal{H})$ . We denote by  $\mathcal{W}_{\mathcal{H}}$  the vanishing set of  $\mathbf{w}_{\mathcal{H}}$ . By Lemma 5.5.1,  $\mathcal{W}_{\mathcal{H}}$  is a proper Zariski closed subset of  $\mathbb{C}^t$ . Our algorithm relies on the following proposition.

**Proposition 5.5.2.** *Assume that Assumption (5.A) holds and the ideal  $\langle \mathbf{f} \rangle$  is radical. Let  $\mathcal{S}$  be a connected component of the semi-algebraic set  $\mathbb{R}^t \setminus (\mathcal{W}_{\infty} \cup \mathcal{W}_{\mathcal{H}})$ . The number of real solutions of  $\mathbf{f}(\eta, \cdot)$  is invariant when  $\eta$  varies over  $\mathcal{S}$ .*

*Proof.* By Lemma 5.4.10,  $\mathcal{W}_{\infty}$  contains the following sets:

- The non-proper points of the restriction of  $\pi$  to  $\mathcal{V}$ .
- The point  $\eta \in \mathbb{C}^t$  such that the fiber  $\pi^{-1}(\eta) \cap \mathcal{V}$  is infinite.
- The image by  $\pi$  of the irreducible components of  $\mathcal{V}$  whose dimensions are smaller than  $t$ .

Now let  $\Delta := \text{jac}(\mathbf{f}, \mathbf{x})$  be the Jacobian matrix of  $\mathbf{f}$  with respect to the variables  $\mathbf{x}$ . The ideal generated by the  $n \times n$ -minors of  $\Delta$  is denoted by  $I_\Delta$ . We consider the algebraic set  $K(\pi, \mathcal{V}) \subset \mathbb{C}^{n+t}$  defined by the ideal  $\langle \mathbf{f} \rangle + I_\Delta$ . Note that, since  $\mathbf{f}$  is radical, by Jacobian criterion,  $K(\pi, \mathcal{V})$  contains the singularities on the irreducible components of  $\mathcal{V}$  of dimension  $t$  and the critical points of the restriction of  $\pi$  to those components.

By Proposition 5.4.8, for  $\eta \in \mathbb{C}^t \setminus \mathcal{W}_\infty$ ,  $\langle \mathbf{f} \rangle$  is a zero-dimensional ideal and the quotient ring  $\mathbb{C}[\mathbf{x}]/\langle \mathbf{f}(\eta, \cdot) \rangle$  has dimension  $\delta$ . Moreover, if  $\eta \in \mathbb{C}^t \setminus (\mathcal{W}_\infty \cup \mathcal{W}_\mathcal{H})$ , the system  $\mathbf{f}(\eta, \cdot)$  has  $\delta$  distinct complex solutions as the rank of  $\mathcal{H}(\eta)$  is  $\delta$ . Therefore, every complex root of  $\mathbf{f}(\eta, \cdot)$  is of multiplicity one (we use the definition of multiplicity given by Proposition 3.3.3).

Now we prove that, for such a point  $\eta$ , the fiber  $\pi^{-1}(\eta)$  does not intersect  $K(\pi, \mathcal{V})$ . Assume by contradiction that there exists a point  $(\eta, \chi) \in \mathbb{C}^{t+n}$  lying in  $\pi^{-1}(\eta) \cap K(\pi, \mathcal{V})$ . Note that  $\chi$  is a solution of  $\mathbf{f}(\eta, \cdot)$ , i.e.,  $\mathbf{f}(\eta, \chi) = 0$ .

As  $(\eta, \chi) \in K(\pi, \mathcal{V})$ , then it is contained in  $V(I_\Delta)$ . Hence, as the derivation in  $\Delta$  does not involve  $\mathbf{y}$ ,  $\chi$  cancels all the  $n \times n$ -minors of the Jacobian matrix  $\text{jac}(\mathbf{f}(\eta, \cdot), \mathbf{x})$ . [9, Proposition 4.16] implies that  $\chi$  has multiplicity greater than one. This contradicts the claim that  $\mathbf{f}(\eta, \cdot)$  admits only complex solutions of multiplicity one. Therefore, we conclude that, for  $\eta \in \mathbb{C}^t \setminus (\mathcal{W}_\infty \cup \mathcal{W}_\mathcal{H})$ ,  $\pi^{-1}(\eta)$  does not intersect  $K(\pi, \mathcal{V})$ .

So, using what we prove above and Lemma 5.4.10, we deduce that, for  $\eta \in \mathbb{R}^t \setminus (\mathcal{W}_\infty \cup \mathcal{W}_\mathcal{H})$ , there exists an open neighborhood  $O_\eta$  of  $\eta$  for the Euclidean topology such that  $\pi^{-1}(O_\eta)$  does not intersect  $K(\pi, \mathcal{V}) \cup \pi^{-1}(\mathcal{W}_\infty)$ . Hence, by Thom's isotopy lemma [47], the restriction of  $\pi$  to  $\mathcal{V} \cap \mathbb{R}^{n+t}$  realizes a locally trivial fibration over  $\mathbb{R}^t \setminus (\mathcal{W}_\infty \cup \mathcal{W}_\mathcal{H})$ . So, for any connected component  $\mathcal{C}$  of  $\mathbb{R}^t \setminus (\mathcal{W}_\infty \cup \mathcal{W}_\mathcal{H})$  and any  $\eta \in \mathcal{C}$ , we have that  $\pi^{-1}(\mathcal{C}) \cap \mathcal{V} \cap \mathbb{R}^{t+n}$  is homeomorphic to  $\mathcal{C} \times (\pi^{-1}(\eta) \cap \mathcal{V} \cap \mathbb{R}^{t+n})$ .

As a consequence, the number of distinct real solutions of  $\mathbf{f}(\eta, \cdot)$  is invariant when  $\eta$  varies over each connected component of  $\mathbb{R}^t \setminus (\mathcal{W}_\infty \cup \mathcal{W}_\mathcal{H})$ .  $\square$

To describe Algorithm 5.3, we need to introduce the following subroutines:

- **CleanFactors** which takes as input a polynomial  $p \in \mathbb{Q}[\mathbf{y}, \mathbf{x}]$  and the polynomial  $w_\infty$ . It computes the square-free part of  $p$  with all the common factors with  $w_\infty$  removed.
- **Signature** which takes as input a symmetric matrix with entries in  $\mathbb{Q}$  and evaluates its signature.
- **SamplePoints** which takes as input a set of polynomials  $g_1, \dots, g_s \in \mathbb{Q}[\mathbf{y}]$  and computes a finite subset  $\mathcal{R}$  of  $\mathbb{Q}^t$  that intersects every connected component of the semi-algebraic set defined by  $\bigwedge_{i=1}^s g_i \neq 0$ . An explicit description of **SamplePoints** is given in the proof of Theorem 5.2.1 in Section 5.2.

The pseudo-code of Algorithm 5.3 is below. Its proof of correctness follows immediately from Proposition 5.5.2 and Corollary 5.4.9.

---

**Algorithm 5.3:** Weak-RRC-Hermite

---

**Input:** A polynomial sequence  $\mathbf{f} \in \mathbb{Q}[\mathbf{y}][x]$  such that  $\langle \mathbf{f} \rangle$  is radical and Assumptions (5.A) holds.

**Output:** A set of sample points and the corresponding numbers of real solutions solving the weak version of Problem (RRC)

- 1  $[\mathcal{H}, \mathbf{w}_\infty] \leftarrow \text{DRL-HermiteMatrix}(\mathbf{f})$
- 2  $\mathbf{w}_\mathcal{H} \leftarrow \text{CleanFactors}(\text{numer}(\det(\mathcal{H})), \mathbf{w}_\infty)$
- 3  $L \leftarrow \text{SamplePoints}(\mathbf{w}_\mathcal{H} \neq 0 \wedge \mathbf{w}_\infty \neq 0)$
- 4 **for**  $\eta \in L$  **do**
- 5      $r_\eta \leftarrow \text{Signature}(\mathcal{H}(\eta))$
- 6 **end**
- 7 **return**  $\{(\eta, r_\eta) \mid \eta \in L\}$

---

**Example 5.5.3.** We continue with the system in Example 5.4.3. The determinant of its parametric Hermite matrix is

$$\mathbf{w}_\mathcal{H} = 16y_1(-y_2^6 - 3y_2^4y_3^2 - 3y_2^2y_3^4 - y_3^6 + 3y_1y_2^4 - 21y_1y_2^2y_3^2 + 3y_1y_3^4 - 3y_1^2y_2^2 - 3y_1^2y_3^2 + y_1^3).$$

We notice that  $\mathbf{w}_\mathcal{H}$  coincides exactly with the output returned by the procedure `DISCRIMINANTVARIETY` of Maple's package `ROOTFINDING[PARAMETRIC]` that computes a discriminant variety [134].

Computing at least one point per connected component of the semi-algebraic set  $\mathbb{R}^3 \setminus V(\mathbf{w}_\mathcal{H})$  using `RAGlib` gives us 28 points. We evaluate the signatures of  $\mathcal{H}$  specialized at those points and find that the input system can have 0, 2 or 4 distinct real solutions when the parameters vary.

**Remark 5.5.4.** As we have seen, Algorithm 5.3 obtains a polynomial which serves similarly as discriminant varieties [134] or border polynomials [203] through computing the determinant of parametric Hermite matrices. The two latter strategies rely on algebraic elimination based on Gröbner bases to compute the projection of  $\text{crit}(\pi, \mathcal{V})$  on the  $\mathbf{y}$ -space. Since it is well-known that the computation of such a Gröbner basis could be expensive, our algorithm has a chance to be more practical. In Section 5.7, we provide experimental results to support this claim.

**Remark 5.5.5.** It is worth noticing that, even though the design of Algorithm 5.3 employs the grevlex monomial ordering where  $x_1 \succ \dots \succ x_n$ , we can replace it by any grevlex ordering with another lexicographical order among the  $x$ 's. For instance, we can use the monomial ordering `grevlex`( $x_n \succ \dots \succ x_1$ ). While the theoretical claims hold for both of these orderings, their practical behaviors could be different. We demonstrate this remark in Example 5.5.6 below.

**Example 5.5.6.** We consider the polynomial sequence  $(f_1, f_2, f_3) \subset \mathbb{Q}[y_1, y_2, y_3][x_1, x_2, x_3]$

$$\begin{aligned} f_1 &= x_1x_2 - x_3, \\ f_2 &= x_1^3 + 4x_1^2x_3 + 2x_2^3 - x_2^2x_3 + x_2x_3^2 - 2x_3^3 + 3x_1^2 - x_1x_3 - 3x_2^2 - 3x_3^2 - x_2 + 4x_3 + 4, \\ f_3 &= y_3x_1x_2 + y_1x_1 + y_2x_2 + 1. \end{aligned}$$

By computing the reduced Gröbner basis of the ideal generated by  $f_1, f_2, f_3$  with respect to the ordering  $\text{grevlex}(x_1 \succ x_2 \succ x_3) \succ \text{grevlex}(y_1 \succ y_2 \succ y_3)$ , one notes that this system above does not satisfy Assumption (5.C). Hence, the algebraic set  $\mathcal{W}_\infty$  defining the locus over which our parametric Hermite matrix does not well specialize is non-empty.

The polynomials  $w_\infty$  and  $w_{\mathcal{H}}$  computed in Algorithm 5.3 with respect to the monomial ordering  $\text{grevlex}(x_1 \succ x_2 \succ x_3)$  have respectively the degrees 13 and 18.

On the other hand, using the monomial ordering  $\text{grevlex}(x_3 \succ x_2 \succ x_1)$  in Algorithm 5.3, one obtains a polynomial  $\tilde{w}_\infty$  of degree 7 and the same polynomial  $w_{\mathcal{H}}$  as above.

Therefore, the degree of the input given to the subroutine `SamplePoints` is reduced by using the second ordering (25 compared with 31). In practice, this choice of ordering accelerates significantly the computation of sample points.

## 5.5.2 Computing semi-algebraic formulas

By Corollary 5.4.9, the number of real roots of the system  $\mathbf{f}(\eta, \cdot)$  for a given point  $\eta \in \mathbb{R}^t \setminus \mathcal{W}_\infty$  can be obtained by evaluating the signature of the parametric Hermite matrix  $\mathcal{H}$ . We recall that the signature of a matrix can be deduced from the sign pattern of its leading principal minors. More precisely, we recall the following criterion, introduced in [191] and [114] (see [77] for a summary of these works).

**Lemma 5.5.7.** [77, Theorem 2.3.6] *Let  $S$  be a  $\delta \times \delta$  symmetric matrix in  $\mathbb{R}^{\delta \times \delta}$  and, for  $1 \leq i \leq \delta$ ,  $S_i$  be the  $i$ -th leading principal minor of  $S$ , i.e., the determinant of the sub-matrix formed by the first  $i$  rows and  $i$  columns of  $S$ . By convention, we denote  $S_0 = 1$ .*

*We assume that  $S_i \neq 0$  for  $0 \leq i \leq \delta$ . Let  $k$  be the number of sign variations between  $S_i$  and  $S_{i+1}$ . Then, the numbers of positive and negative eigenvalues of  $S$  are respectively  $\delta - k$  and  $k$ . Thus, the signature of  $S$  is  $\delta - 2k$ .*

This criterion leads us to the following idea. Assume that none of the leading principal minors of  $\mathcal{H}$  is identically zero. We consider the semi-algebraic subset of  $\mathbb{R}^t$  defined by the non-vanishing of those leading principal minors. Over a connected component  $\mathcal{S}'$  of this semi-algebraic set, each leading principal minor is not zero and its sign is invariant. As a consequence, by Lemma 5.5.7 and Corollary 5.4.9, the number of distinct real roots of  $\mathbf{f}(\eta, \cdot)$  when  $\eta$  varies over  $\mathcal{S}' \setminus \mathcal{W}_\infty$  is invariant.

However, this approach does not apply directly if one of the leading principal minors of  $\mathcal{H}$  is identically zero. We bypass this obstacle by picking randomly an invertible matrix  $A \in \text{GL}(\delta, \mathbb{Q})$  and working with the matrix  $\mathcal{H}_Q := Q^T \cdot \mathcal{H} \cdot Q$ . The lemma below states that, with a generic matrix  $Q$ , all of the leading principal minors of  $\mathcal{H}_Q$  are not identically zero.

**Lemma 5.5.8.** *There exists a Zariski dense subset  $\mathcal{Q}$  of  $\mathrm{GL}(\delta, \mathbb{Q})$  such that for  $Q \in \mathcal{Q}$ , all of the leading principal minors of  $\mathcal{H}_Q := Q^T \cdot \mathcal{H} \cdot Q$  are not identically zero.*

*Proof.* For  $1 \leq r \leq \delta$ , we denote by  $\mathfrak{M}_r$  the set of all  $r \times r$  minors of  $\mathcal{H}$ .

Let  $\eta \in \mathbb{Q}^t \setminus \mathcal{W}_\infty \cup \mathcal{W}_\mathcal{H}$ . We have that  $\mathcal{H}(\eta)$  is a full rank matrix in  $\mathbb{Q}^{\delta \times \delta}$  and, for  $Q \in \mathrm{GL}(\delta, \mathbb{R})$ ,  $\mathcal{H}_Q(\eta) = Q^T \cdot \mathcal{H}(\eta) \cdot Q$ .

We prove that there exists a Zariski dense subset  $\mathcal{Q}$  of  $\mathrm{GL}(\delta, \mathbb{Q})$  such that, for  $Q \in \mathcal{Q}$ , all of the leading principal minors of  $\mathcal{H}_Q(\eta)$  are not zero. Then, as an immediate consequence, all the leading principal minors of  $\mathcal{H}_Q$  are not identically zero.

We consider the matrix  $Q = (q_{i,j})_{1 \leq i,j \leq \delta}$  where  $\mathbf{q} = (q_{i,j})$  are new variables. Then, the  $r$ -th leading principal minor  $M_r(\mathbf{q})$  of  $Q^T \cdot \mathcal{H}(\eta) \cdot Q$  can be written as

$$M_r(\mathbf{q}) = \sum_{\mathfrak{m} \in \mathfrak{M}_r} q_{\mathfrak{m}} \cdot \mathfrak{m}(\eta),$$

where the  $q_{\mathfrak{m}}$ 's are elements of  $\mathbb{Q}[\mathbf{q}]$ .

As  $\mathcal{H}(\eta)$  is a full rank symmetric matrix by assumption, there exists a matrix  $Q \in \mathrm{GL}(\delta, \mathbb{R})$  such that  $Q^T \cdot \mathcal{H}(\eta) \cdot Q$  is a diagonal matrix with no zero on its diagonal. Hence, the evaluation of  $\mathbf{q}$  at the entries of  $Q$  gives  $M_r(\mathbf{q})$  a non-zero value. As a consequence,  $M_r(\mathbf{q})$  is not identically zero.

Let  $\mathcal{Q}_r$  be the non-empty Zariski open subset of  $\mathrm{GL}(\delta, \mathbb{Q})$  defined by  $M_r(\mathbf{q}) \neq 0$ . Then, the set of the matrices  $Q \in \mathcal{Q}_r$  such that the  $r \times r$  leading principal minor of  $Q^T \cdot \mathcal{H}(\eta) \cdot Q$  is not zero.

Taking  $\mathcal{Q}$  as the intersection of  $\mathcal{Q}_r$  for  $1 \leq r \leq \delta$ , then, for  $Q \in \mathcal{Q}$ , none of the leading principal minors of  $Q^T \cdot \mathcal{H}(\eta) \cdot Q$  equals zero. Consequently, each leading principal minor of  $Q^T \cdot \mathcal{H} \cdot Q$  is not identically zero.  $\square$

Our algorithm (Algorithm 5.4) for solving Problem (RRC) through parametric Hermite matrices is described below. As it depends on the random choice of the matrix  $Q$ , Algorithm 5.4 is probabilistic. Note that one can easily detect the cancellation of the leading principal minors for each choice of  $Q$ .

---

**Algorithm 5.4: RRC-Hermite**

---

**Input:** A polynomial sequence  $\mathbf{f} \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  such that the ideal  $\langle \mathbf{f} \rangle$  is radical and  $\mathbf{f}$  satisfies Assumption (5.A)

**Output:** The descriptions of a collection of semi-algebraic sets  $\mathcal{S}_i$  solving Problem (RRC)

```
1  $\mathcal{H}, \mathbf{w}_\infty \leftarrow \text{DRL-HermiteMatrix}(\mathbf{f})$ 
2 Choose randomly a matrix  $Q$  in  $\text{GL}(\delta, \mathbb{Q})$ 
3  $\mathcal{H}_Q \leftarrow Q^T \cdot \mathcal{H} \cdot Q$ 
4  $(M_1, \dots, M_\delta) \leftarrow \text{LeadingPrincipalMinors}(\mathcal{H}_Q)$ 
5  $(m_1, \dots, m_\delta) \leftarrow \text{Numerators}(M_1, \dots, M_\delta)$ 
6  $L \leftarrow \text{SamplePoints}((\wedge_{i=1}^\delta m_i \neq 0) \wedge \mathbf{w}_\infty \neq 0)$ 
7 for  $\eta \in L$  do
8    $r_\eta \leftarrow \text{Signature}(\mathcal{H}(\eta))$ 
9 end
10  $\mathbf{w}_f \leftarrow \mathbf{w}_\infty \cdot m_1 \cdot \dots \cdot m_\delta$ 
11 return  $\{(\text{sign}(M_1(\eta)), \dots, M_\delta(\eta)), \eta, r_\eta \mid \eta \in L\}$  and  $\mathbf{w}_f$ 
```

---

**Proposition 5.5.9.** *Assume that  $\mathbf{f}$  satisfies Assumptions (5.A) and that the ideal  $\langle \mathbf{f} \rangle$  is radical. Let  $Q$  be a matrix in  $\text{GL}(\delta, \mathbb{Q})$  such that all of the leading principal minors  $M_1, \dots, M_\delta$  of  $\mathcal{H}_Q := Q^T \cdot \mathcal{H} \cdot Q$  are not identically zero. Then, Algorithm 5.4 computes correctly a solution for Problem (RRC).*

*Proof.* Note that for  $\eta \in \mathbb{R}^t \setminus \mathcal{W}_\infty$ , we have that  $\mathcal{H}_Q(\eta) = Q^T \cdot \mathcal{H}(\eta) \cdot Q$ . Therefore, the signature of  $\mathcal{H}(\eta)$  equals to the signature of  $\mathcal{H}_Q(\eta)$ .

Let  $M_1, \dots, M_\delta$  be the leading principal minors of  $\mathcal{H}_Q$  and  $\mathcal{S}$  be the algebraic set defined by  $\wedge_{i=1}^\delta M_i \neq 0$ . Over each connected component  $\mathcal{S}'$  of  $\mathcal{S}$ , the sign of each  $M_i$  is invariant and not zero. Therefore, by Lemma 5.5.7, the signature of  $\mathcal{H}_Q(\eta)$ , and therefore of  $\mathcal{H}(\eta)$ , is invariant when  $\eta$  varies over  $\mathcal{S}' \setminus \mathcal{W}_\infty$ . As a consequence, by Corollary 5.4.9, the number of distinct real roots of  $\mathbf{f}(\eta, \cdot)$  is also invariant when  $\eta$  varies over  $\mathcal{S}' \setminus \mathcal{W}_\infty$ . We finish the proof of correctness of Algorithm 5.4.  $\square$

**Example 5.5.10.** *From the parametric Hermite matrix  $\mathcal{H}$  computed in Example 5.4.3, we obtain the sequence of leading principal minors below:*

$$\begin{aligned} M_1 &= 4, \\ M_2 &= 4(-2y_2^2 + y_3^2 + 2y_1), \\ M_3 &= 8(-y_2^4 - 2y_2^2y_3^2 - y_3^4 - y_1y_2^2 - y_1y_3^2 + 2y_1^2), \\ M_4 &= 16y_1(-y_2^6 - 3y_2^4y_3^2 - 3y_2^2y_3^4 - y_3^6 + 3y_1y_2^4 - 21y_1y_2^2y_3^2 + 3y_1y_3^4 - 3y_1^2y_2^2 - 3y_1^2y_3^2 + y_1^3). \end{aligned}$$

Since  $M_1$  is constant, we compute at least one point per connected component of the semi-algebraic set defined by

$$M_2 \neq 0 \wedge M_3 \neq 0 \wedge M_4 \neq 0.$$

The computation using RAGlib outputs a set of 48 sample points and finds the following realizable sign conditions of  $(M_2, M_3, M_4)$ :

$$[-1, 1, 1], [-1, -1, 1], [1, -1, -1], [-1, -1, -1], [1, 1, -1].$$

By evaluating the signature of  $\mathcal{H}$  at each of those sample points, we deduce the semi-algebraic formulas corresponding to every possible number of real solutions

$$\begin{aligned} 0 \text{ real root} &\rightarrow (M_2 < 0 \wedge M_3 > 0 \wedge M_4 > 0) \vee (M_2 < 0 \wedge M_3 < 0 \wedge M_4 > 0), \\ 2 \text{ real roots} &\rightarrow (M_2 > 0 \wedge M_3 < 0 \wedge M_4 < 0) \vee (M_2 < 0 \wedge M_3 < 0 \wedge M_4 < 0) \\ &\quad \vee (M_2 > 0 \wedge M_3 > 0 \wedge M_4 < 0), \\ 4 \text{ real roots} &\rightarrow (M_2 > 0 \wedge M_3 > 0 \wedge M_4 > 0). \end{aligned}$$

We recall that the semi-algebraic formulas obtained in Example 5.3.3 involve the subresultant coefficients  $s_2$ ,  $s_3$  and  $s_4$  of degree 2, 6 and 11 respectively. Whereas, the degrees of the minors  $M_2$ ,  $M_3$  and  $M_4$  that we obtain from the parametric Hermite matrix are only 2, 4 and 7.

## 5.6 Complexity analysis

This section is devoted to the complexity analysis of Algorithms 5.2, 5.3 and 5.4. Under some genericity assumptions, we provide degree bounds for entries of parametric Hermite matrices constructed using grevlex ordering and polynomials involving in our output semi-algebraic formulas, which are actually minors of those matrices. Using these bounds, we prove that our algorithms run within  $D^{O(nt)}$  arithmetic operations in  $\mathbb{Q}$  on generic inputs.

### 5.6.1 Degree bounds of parametric Hermite matrices on generic input

We consider an affine regular sequence  $\mathbf{f} = (f_1, \dots, f_n) \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  according to the variables  $\mathbf{x}$ , i.e., the homogeneous components of largest degree in  $\mathbf{x}$  of the  $f_i$ 's form a homogeneous regular sequence. Additionally, we require that  $\mathbf{f}$  satisfies Assumptions (5.A) and (5.C).

Let  $D$  be the highest value among the total degrees of the  $f_i$ 's. Since the homogeneous regular sequences are generic among the homogeneous polynomial sequences (Proposition 2.7.8), the same property of genericity holds for affine regular sequences (thanks to the definition we use).

As in previous sections,  $\mathcal{G}$  denotes the reduced Gröbner basis of  $\langle \mathbf{f} \rangle$  with respect to the ordering  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$ . Let  $\delta$  be the dimension of the  $\mathbb{K}$ -vector space  $\mathbb{K}[\mathbf{x}]/\langle \mathbf{f} \rangle_{\mathbb{K}}$  where  $\mathbb{K} = \mathbb{Q}(\mathbf{y})$ . By Bézout's inequality,  $\delta \leq D^n$ . We derive from  $\mathcal{G}$  a basis  $\mathcal{B} = \{b_1, \dots, b_\delta\}$  of  $\mathbb{K}[\mathbf{x}]/\langle \mathbf{f} \rangle_{\mathbb{K}}$  consisting of monomials in the variables  $\mathbf{x}$ . Finally, the parametric Hermite matrix of  $\mathbf{f}$  with respect to  $\mathcal{B}$  is denoted by  $\mathcal{H} = (h_{i,j})_{1 \leq i,j \leq \delta}$ .

For a polynomial  $p \in \mathbb{Q}[\mathbf{y}, \mathbf{x}]$ , we denote by  $\deg(p)$  the total degree of  $p$  in  $(\mathbf{y}, \mathbf{x})$  and  $\deg_{\mathbf{x}}(p)$  the partial degree of  $p$  in the variables  $\mathbf{x}$ .

As Assumption (5.C) holds, by Lemma 5.4.5, the entries of the parametric Hermite matrix  $\mathcal{H}$  associated to  $\mathbf{f}$  with respect to the basis  $\mathcal{B}$  are elements of  $\mathbb{Q}[\mathbf{y}]$ . To establish a degree bound on the entries of  $\mathcal{H}$ , we need to introduce the following assumption.

**Assumption 5.D.** For any  $g \in \mathcal{G}$ , we have that  $\deg(g) = \deg_{\mathbf{x}}(g)$ .

Proposition 5.6.1 below states that Assumption (5.D) is generic. Its direct consequence is a proof for Proposition 5.4.6.

**Proposition 5.6.1.** Let  $\mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}$  be the set of polynomials in  $\mathbb{C}[\mathbf{x}, \mathbf{y}]$  having total degree bounded by  $D$ . There exists a non-empty Zariski open subset  $\mathcal{F}_D$  of  $\mathbb{C}[\mathbf{x}, \mathbf{y}]_D^n$  such that Assumption (5.D) holds for  $\mathbf{f} \in \mathcal{F}_D \cap \mathbb{Q}[\mathbf{x}, \mathbf{y}]^n$ .

Consequently, for  $\mathbf{f} \in \mathcal{F}_D \cap \mathbb{Q}[\mathbf{x}, \mathbf{y}]^n$ ,  $\mathbf{f}$  satisfies Assumption (5.C).

*Proof.* Let  $y_{t+1}$  be a new indeterminate. For any polynomial  $p \in \mathbb{Q}[\mathbf{x}, \mathbf{y}]$ , we consider the homogenized polynomial  $p_h \in \mathbb{Q}[\mathbf{x}, \mathbf{y}, y_{t+1}]$  of  $p$  defined as follows:

$$p_h = y_{t+1}^{\deg(p)} p \left( \frac{x_1}{y_{t+1}}, \dots, \frac{x_n}{y_{t+1}}, \frac{y_1}{y_{t+1}}, \dots, \frac{y_t}{y_{t+1}} \right).$$

Let  ${}^h\mathbb{C}[\mathbf{x}, \mathbf{y}, y_{t+1}]_D$  be the set of homogeneous polynomials in  $\mathbb{C}[\mathbf{x}, \mathbf{y}, y_{t+1}]$  whose degrees are exactly  $D$ . By Proposition 2.7.8, there exists a non-empty Zariski subset  ${}^h\mathcal{F}_D$  of  ${}^h\mathbb{C}[\mathbf{x}, \mathbf{y}, y_{t+1}]_D^n$  such that the variables  $\mathbf{x}$  is in Noether position with respect to  ${}^h\mathbf{f}$  for every  ${}^h\mathbf{f} \in {}^h\mathcal{F}_D$ .

For  ${}^h\mathbf{f} \in {}^h\mathcal{F}_D$ , let  ${}^hG$  be the reduced Gröbner basis of  ${}^h\mathbf{f}$  with respect to the grevlex ordering  $\text{grevlex}(\mathbf{x} \succ \mathbf{y} \succ y_{t+1})$ . By [7, Proposition 7], if the variables  $\mathbf{x}$  is in Noether position with respect to  ${}^h\mathbf{f}$ , then the leading monomials appearing in  ${}^hG$  depend only on  $\mathbf{x}$ .

Let  $\mathbf{f}$  and  $G$  be the image of  ${}^h\mathbf{f}$  and  ${}^hG$  by substituting  $y_{t+1} = 1$ . We show that  $G$  is a Gröbner basis of  $\mathbf{f}$  with respect to the ordering  $\text{grevlex}(\mathbf{x} \succ \mathbf{y})$ .

Since  ${}^hG$  generates  $\langle {}^h\mathbf{f} \rangle$ ,  $G$  is a generating set of  $\langle \mathbf{f} \rangle$ . As the leading monomials of elements in  ${}^hG$  do not depend on  $y_{t+1}$ , the substitution  $y_{t+1} = 1$  does not affect these leading monomials.

For a polynomial  $p \in \langle \mathbf{f} \rangle \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$ , then  $p$  writes

$$p = \sum_{i=1}^n c_i \cdot f_i,$$

where the  $c_i$ 's lie in  $\mathbb{Q}[\mathbf{x}, \mathbf{y}]$ . By homogenizing each polynomial  $c_i \cdot f_i$  on the right hand side, we obtain a homogeneous polynomial  ${}^hP \in \langle {}^h\mathbf{f} \rangle$ :

$${}^hP = \sum_{i=1}^n y_{t+1}^{\bullet} \cdot {}^h c_i \cdot {}^h f_i$$

where  $y_{t+1}^{\bullet}$  means some suitable power of  $y_{t+1}$ .

Note that  ${}^hP$  is not necessarily the homogenization  ${}^h p$  of  $p$  but only the product of  ${}^h p$  with a power of  $y_{t+1}$ . Then, there exists a polynomial  ${}^h g \in {}^hG$  such that the leading monomial of  ${}^h g$  divides the leading monomial of  ${}^h P$ . Since the leading monomial of  ${}^h g$  depends only on  $\mathbf{x}$ , it also divides the leading monomial of  ${}^h p$ , which is the leading monomial of  $p$ . So, the leading monomial of the image of  ${}^h g$  in  $G$  divides the leading monomial of  $p$ . We conclude that  $G$  is a Gröbner basis

of  $\mathbf{f}$  with respect to the ordering  $\text{grevlex}(\mathbf{x} \succ \mathbf{y})$  and the set of leading monomials in  $G$  depends only on the variables  $\mathbf{x}$ .

Let  $\mathcal{F}_D$  be the subset of  $\mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^n$  such that for every  $\mathbf{f} \in \mathcal{F}_D$ , its homogenization  ${}^h\mathbf{f}$  is contained in  ${}^h\mathcal{F}_D$ .

Since the two spaces  ${}^h\mathbb{C}[\mathbf{x}, \mathbf{y}, y_{t+1}]_D^n$  and  $\mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^n$  are both isomorphic to the affine space  $\mathbb{C}^{\binom{d+n+t}{n+t} \times n}$  (by considering each monomial coefficient as a coordinate),  $\mathcal{F}_D$  is also a non-empty Zariski open subset of  $\mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^n$ .

Assume now that the polynomial sequence  $\mathbf{f}$  belongs to  $\mathcal{F}_D$ . We consider the two monomial orderings over  $\mathbb{Q}[\mathbf{x}, \mathbf{y}]$  below:

- The elimination ordering  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$  is abbreviated by  $O_1$ . The leading monomial of  $p \in \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  with respect to  $O_1$  is denoted by  $\text{lm}_1(p)$ . The reduced Gröbner basis of  $\mathbf{f}$  with respect to  $O_1$  is  $\mathcal{G}$ .
- The grevlex ordering  $\text{grevlex}(\mathbf{x} \succ \mathbf{y})$  is abbreviated by  $O_2$ . The leading monomial of  $p \in \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  with respect to  $O_2$  is denoted by  $\text{lm}_2(p)$ . The reduced Gröbner basis of  $\mathbf{f}$  with respect to  $O_2$  is denoted by  $\mathcal{G}_2$ .

As proven above, the set  $\{\text{lm}_2(g_2) \mid g_2 \in \mathcal{G}_2\}$  does not depend on  $\mathbf{y}$ . With this property, we will show, for any  $g_2 \in \mathcal{G}_2$ , there exists a polynomial  $g \in \mathcal{G}$  such that  $\text{lm}_1(g)$  divides  $\text{lm}_2(g_2)$ .

By definition,  $\text{lm}_2(g_2)$  is greater than any other monomial of  $g_2$  with respect to the ordering  $O_2$ . Since  $\text{lm}_2(g_2)$  depends only on the variables  $\mathbf{x}$ , it is then greater than any monomial of  $g_2$  with respect to the ordering  $O_1$ . Hence,  $\text{lm}_2(g_2)$  is also  $\text{lm}_1(g_2)$ . Consequently, since  $\mathcal{G}$  is a Gröbner basis of  $\mathbf{f}$  with respect to  $O_1$ , there exists a polynomial  $g \in \mathcal{G}$  such that  $\text{lm}_1(g)$  divides  $\text{lm}_1(g_2) = \text{lm}_2(g_2)$ .

Next, we prove that for every  $g \in \mathcal{G}$ ,  $\text{lm}_1(g)$  is also  $\text{lm}_2(g)$ . For this, we rely on the fact that  $\mathcal{G}$  is reduced. Assume by contradiction that there exists a polynomial  $g \in \mathcal{G}$  such that  $\text{lm}_1(g) \neq \text{lm}_2(g)$ . Thus,  $\text{lm}_2(g)$  must contain both  $\mathbf{x}$  and  $\mathbf{y}$ . Let  $t_{\mathbf{x}}$  be the part in only variables  $\mathbf{x}$  of  $\text{lm}_2(g)$ . Note that  $\text{lm}_1(g)$  is greater than  $t_{\mathbf{x}}$  with respect to  $O_1$ . There exists an element  $g_2 \in \mathcal{G}_2$  such that  $\text{lm}_2(g_2)$  divides  $\text{lm}_2(g)$ . Since  $\text{lm}_2(g_2)$  depends only on the variables  $\mathbf{x}$ , we have that  $\text{lm}_2(g_2)$  divides  $t_{\mathbf{x}}$ . Then, by what we proved above, there exists  $g' \in \mathcal{G}$  such that  $\text{lm}_1(g')$  divides  $\text{lm}_2(g_2)$ , so  $\text{lm}_1(g')$  divides  $t_{\mathbf{x}}$ . This implies that  $\mathcal{G}$  is not reduced, which contradicts the definition of  $\mathcal{G}$ .

So,  $\text{lm}_1(g) = \text{lm}_2(g)$  for every  $g \in \mathcal{G}$  and, consequently,  $\deg(g) = \deg_{\mathbf{x}}(g)$ .

We conclude that there exists a non-empty Zariski open subset  $\mathcal{F}_D$  (as above) of  $\mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^n$  such that Assumption (5.D) holds for every  $\mathbf{f} \in \mathcal{F}_D \cap \mathbb{Q}[\mathbf{x}, \mathbf{y}]^n$ .

Additionally, one easily notices that Assumption (5.D) implies Assumption (5.C). As a consequence,  $\mathbf{f}$  also satisfies Assumption (5.C) for any  $\mathbf{f} \in \mathcal{F}_D \cap \mathbb{Q}[\mathbf{x}, \mathbf{y}]^n$ .  $\square$

Recall that, when Assumption (5.C) holds, by Lemma 5.4.5, the trace of any multiplication map  $\mathcal{L}_p$  is a polynomial in  $\mathbb{Q}[\mathbf{y}]$  where  $p \in \mathbb{Q}[\mathbf{y}][\mathbf{x}]$ . We now estimate the degree of  $\text{trace}(\mathcal{L}_p)$ . Since the map  $p \mapsto \text{trace}(\mathcal{L}_p)$  is linear, it is sufficient to consider  $p$  as a monomial in the variables  $\mathbf{x}$ .

**Proposition 5.6.2.** *Assume that Assumption (5.D) holds. Then, for any monomial  $m$  in the variables  $\mathbf{x}$ , the degree in  $\mathbf{y}$  of  $\text{trace}(\mathcal{L}_m)$  is bounded by  $\deg(m)$ . As a consequence, the total degree of the entry  $h_{i,j} = \text{trace}(\mathcal{L}_{b_i \cdot b_j})$  of  $\mathcal{H}$  is at most the sum of the total degrees of  $b_i$  and  $b_j$ , i.e.,*

$$\deg(h_{i,j}) \leq \deg(b_i) + \deg(b_j).$$

*Proof.* Let  $m$  be a monomial in  $\mathbb{Q}[\mathbf{x}]$ . The multiplication matrix  $\mathcal{L}_m$  is built as follows. For  $1 \leq i \leq \delta$ , the normal form of  $b_i \cdot m$  as a polynomial in  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$  writes

$$\text{NF}_{\mathcal{G}}(b_i \cdot m) = \sum_{j=1}^{\delta} c_{i,j} \cdot b_j.$$

Note that this normal form is the remainder of the successive divisions of  $b_i \cdot m$  by polynomials in  $\mathcal{G}$ . As Assumption (5.D) holds, Assumption (5.C) also holds. Therefore, those divisions do not introduce any denominator. So, every term appearing during these normal form reductions are polynomials in  $\mathbb{Q}[\mathbf{y}][\mathbf{x}]$ .

Let  $p \in \mathbb{Q}[\mathbf{y}][\mathbf{x}]$ . For any  $g \in \mathcal{G}$ , by Assumption (5.D), the total degree in  $(\mathbf{y}, \mathbf{x})$  of every term of  $g$  is at most the degree of  $\text{lm}_{\mathbf{x}}(g)$ . Thus, a division of  $p$  by  $g$  involves only terms of total degree  $\deg(p)$ . Thus, during the polynomial division of  $p$  to  $\mathcal{G}$ , only terms of degree at most  $\deg(p)$  will appear. Hence the degree of  $\text{NF}_{\mathcal{G}}(p)$  is bounded by  $\deg(p)$ .

Note that  $\text{trace}(\mathcal{L}_m) = \sum_{i=1}^{\delta} c_{i,i}$ . As the degree of  $c_{i,i} \cdot b_i$  is bounded by  $\deg(b_i) + \deg(m)$ , the degree of  $c_{i,i}$  is at most  $\deg(m)$ . Then, we obtain that  $\deg(\text{trace}(\mathcal{L}_m)) \leq \deg(m)$ .

Finally, the degree bound of  $h_{i,j}$  follows immediately:

$$\deg(h_{i,j}) = \deg(\text{trace}(\mathcal{L}_{b_i \cdot b_j})) \leq \deg(b_i \cdot b_j) = \deg(b_i) + \deg(b_j).$$

□

**Lemma 5.6.3.** *Assume that  $\mathbf{f}$  satisfies Assumption (5.D). Then the degree of a minor  $M$  consisting of the rows  $(r_1, \dots, r_{\ell})$  and the columns  $(c_1, \dots, c_{\ell})$  of  $\mathcal{H}$  is bounded by*

$$\sum_{i=1}^{\ell} (\deg(b_{r_i}) + \deg(b_{c_i})).$$

*Particularly, the degree of  $\det(\mathcal{H})$  is bounded by  $2 \sum_{i=1}^{\delta} \deg(b_i)$ .*

*Proof.* We expand the minors  $M$  into terms of the form  $(-1)^{\text{sign}(\sigma)} h_{r_1, \sigma(c_1)} \dots h_{r_{\ell}, \sigma(c_{\ell})}$ , where  $\sigma$  is a permutation of  $\{c_1, \dots, c_{\ell}\}$  and  $\text{sign}(\sigma)$  is its signature. We then bound the degree of each of those terms as follows using Proposition 5.6.2:

$$\begin{aligned} \deg \left( \prod_{i=1}^{\ell} h_{r_i, \sigma(c_i)} \right) &= \sum_{i=1}^{\ell} \deg(h_{r_i, \sigma(c_i)}) \\ &\leq \sum_{i=1}^{\ell} (\deg(b_{r_i}) + \deg(b_{\sigma(c_i)})) = \sum_{i=1}^{\ell} (\deg(b_{r_i}) + \deg(b_{c_i})). \end{aligned}$$

Hence, taking the sum of all those terms, we obtain the inequality:

$$\deg(M_i) \leq \sum_{i=1}^{\ell} (\deg(b_{r_i}) + \deg(b_{c_i})).$$

When  $M$  is taken as the determinant of  $\mathcal{H}$ , then

$$\deg(\det(\mathcal{H})) \leq 2 \sum_{i=1}^{\delta} \deg(b_i).$$

□

Proposition 5.6.2 implies that, when Assumption (5.D) holds, the degree pattern of  $\mathcal{H}$  depends only on the degree of the elements of  $\mathcal{B} = \{b_1, \dots, b_{\delta}\}$ . We rearrange  $\mathcal{B}$  in the increasing order of degree, i.e.,  $\deg(b_i) \leq \deg(b_j)$  for  $1 \leq i < j \leq \delta$ . So,  $b_1 = 1$  and  $\deg(b_1) = 0$ . The degree bounds of the entries of  $\mathcal{H}$  are expressed by the matrix below

$$\begin{bmatrix} 0 & \deg(b_2) & \dots & \deg(b_{\delta}) \\ \deg(b_2) & 2 \deg(b_2) & \dots & \deg(b_{\delta}) + \deg(b_2) \\ \vdots & \vdots & \ddots & \vdots \\ \deg(b_{\delta}) & \deg(b_{\delta}) + \deg(b_2) & \dots & 2 \deg(b_{\delta}) \end{bmatrix}.$$

Moreover, using the regularity of  $\mathbf{f}$ , we are able to establish explicit degree bounds for the elements of  $\mathcal{B}$  and then, for the minors of  $\mathcal{H}$ .

**Lemma 5.6.4.** *Assume that  $\mathbf{f}$  is an affine regular sequence and let  $\mathcal{B}$  be the basis defined as above. Then the highest degree among the elements of  $\mathcal{B}$  is bounded by  $n(D-1) + 1$  and*

$$2 \sum_{i=1}^{\delta} \deg(b_i) \leq n(D-1)D^n.$$

*Proof.* For  $p \in \mathbb{K}[\mathbf{x}]$ , let  ${}^h p \in \mathbb{K}[x_1, \dots, x_{n+1}]$  be the homogenization of  $p$  with respect to the variable  $x_{n+1}$ , i.e.,

$${}^h p = x_{n+1}^{\deg_{\mathbf{x}}(p)} p \left( \frac{x_1}{x_{n+1}}, \dots, \frac{x_n}{x_{n+1}} \right).$$

The dehomogenization map  $\alpha$  is defined as:

$$\begin{aligned} \alpha : \mathbb{K}[x_1, \dots, x_{n+1}] &\rightarrow \mathbb{K}[x_1, \dots, x_n], \\ p(x_1, \dots, x_{n+1}) &\mapsto p(x_1, \dots, x_n, 1). \end{aligned}$$

Also, the homogeneous component of largest degree of  $p$  with respect to the variables  $\mathbf{x}$  is denoted by  ${}^H p$ . Throughout this proof, we use the following notations:

- $I = \langle \mathbf{f} \rangle_{\mathbb{K}}$  and  $\mathcal{G}$  is the reduced Gröbner basis of  $I$  with respect to  $\text{grevlex}(x_1 \succ \cdots \succ x_n)$ .
- ${}^h I = \langle {}^h f_1, \dots, {}^h f_n \rangle_{\mathbb{K}}$  and  ${}^h \mathcal{G}$  is the reduced Gröbner basis of  ${}^h I$  with respect to the ordering  $\text{grevlex}(x_1 \succ \cdots \succ x_{n+1})$ .

The Hilbert series of the homogeneous ideal  ${}^h I$  writes

$$\text{HS}_{{}^h I}(z) = \sum_{r=0}^{\infty} \left( \dim_{\mathbb{K}} \mathbb{K}[\mathbf{x}]_r - \dim_{\mathbb{K}} ({}^h I \cap \mathbb{K}[\mathbf{x}]_r) \right) \cdot z^r,$$

where  $\mathbb{K}[\mathbf{x}]_r = \{p \mid p \in \mathbb{K}[\mathbf{x}] : \deg_{\mathbf{x}}(p) = r\}$ .

Since  $\mathbf{f}$  is an affine regular sequence, by definition,  ${}^H \mathbf{f} = ({}^H f_1, \dots, {}^H f_n)$  forms a homogeneous regular sequence. Equivalently, by [194, Proposition 1.44], the homogeneous polynomial sequence  $({}^h f_1, \dots, {}^h f_n, x_{n+1})$  is regular. Particularly,  $({}^h f_1, \dots, {}^h f_n)$  is a homogeneous regular sequence and, by [155, Theorem 1.5], we obtain

$$\text{HS}_{{}^h I}(z) = \frac{\prod_{i=1}^n (1 - z^{\deg(f_i)})}{(1 - z)^{n+1}} = \frac{\prod_{i=1}^n (1 + \dots + z^{\deg(f_i)-1})}{1 - z}.$$

On the other hand, as  $({}^h f_1, \dots, {}^h f_n, x_{n+1})$  is a homogeneous regular sequence, by [7, Proposition 7], the leading terms of  ${}^h \mathcal{G}$  with respect to  $\text{grevlex}(x_1 \succ \cdots \succ x_{n+1})$  do not depend on the variables  $x_{n+1}$ . Thus, the dehomogenization map  $\alpha$  does not affect the set of leading terms of  ${}^h \mathcal{G}$ . Besides,  $\alpha({}^h \mathcal{G})$  is a Gröbner basis of  $I$  with respect to  $\text{grevlex}(\mathbf{x})$  (see, e.g., the proof of [187, Lemma 27]). Hence, the leading terms of  ${}^h \mathcal{G}$  coincides with the leading terms of  $\mathcal{G}$ .

As a consequence, the set of monomials in  $(x_1, \dots, x_{n+1})$  which are not contained in the initial ideal of  ${}^h I$  with respect to  $\text{grevlex}(x_1 \succ \cdots \succ x_{n+1})$  is exactly

$$\{b \cdot x_{n+1}^j \mid b \in \mathcal{B}, j \in \mathbb{N}\}.$$

Therefore, we have the following equality

$$\dim_{\mathbb{K}} \mathbb{K}[\mathbf{x}]_r - \dim_{\mathbb{K}} ({}^h I \cap \mathbb{K}[\mathbf{x}]_r) = \{b \in \mathcal{B} \mid \deg(b) \leq r\} = \sum_{j=0}^r |\mathcal{B} \cap \mathbb{K}[\mathbf{x}]_j|.$$

Let  $H(z) = \sum_{r=0}^{\infty} |\mathcal{B} \cap \mathbb{K}[\mathbf{x}]_r| \cdot z^r$ . We have that

$$(1 - z) \cdot \text{HS}_{{}^h I}(z) = (1 - z) \sum_{r=0}^{\infty} \sum_{j=0}^r |\mathcal{B} \cap \mathbb{K}[\mathbf{x}]_j| \cdot z^r = \sum_{r=0}^{\infty} |\mathcal{B} \cap \mathbb{K}[\mathbf{x}]_r| \cdot z^r = H(z).$$

Then,

$$H(z) = \prod_{i=1}^n (1 + \dots + z^{\deg(f_i)-1}).$$

As a direct consequence, we obtain the bound

$$\max_{1 \leq i \leq \delta} \deg(b_i) \leq \sum_{i=1}^n \deg(f_i) - n \leq n(D-1).$$

Let  $G_1$  and  $G_2$  be two polynomials in  $\mathbb{Z}[z]$ . We write  $G_1 \leq G_2$  if and only if for any  $r \geq 0$ , the coefficient of  $z^r$  in  $G_2$  is greater than or equal to the one in  $G_1$ .

Since  $\deg(f_i) \leq d$  for every  $1 \leq i \leq n$ , then

$$H(z) = \prod_{i=1}^n (1 + \dots + z^{\deg(f_i)-1}) \leq \prod_{i=1}^n (1 + \dots + z^{D-1}).$$

As a consequence,

$$H'(z) = \sum_{r=1}^{\infty} (r |\mathcal{B} \cap \mathbb{K}[\mathbf{x}]_r|) \cdot z^{r-1} \leq \left( \prod_{i=1}^n (1 + \dots + z^{D-1}) \right)'$$

Expanding  $G'(z)$ , we obtain

$$\begin{aligned} H'(z) &\leq \frac{n(1 + \dots + z^{D-1})^{n-1} (1 + \dots + z^{D-1} - dz^{D-1})}{1-z} \\ &= n(1 + \dots + z^{D-1})^{n-1} \sum_{i=0}^{D-2} \frac{z^i - z^{D-1}}{1-z} \\ &= n(1 + \dots + z^{D-1})^{n-1} \sum_{i=0}^{D-2} z^i (1 + \dots + z^{D-i-2}). \end{aligned}$$

By substituting  $z = 1$  in the above inequality, we obtain

$$H'(1) \leq nD^{n-1} \sum_{i=0}^{D-2} (D-i-1) = \frac{n(D-1)D^n}{2}.$$

Thus, we have that

$$\sum_{i=1}^{\delta} \deg(b_i) = \sum_{r=0}^{\infty} r |\mathcal{B} \cap \mathbb{K}[\mathbf{x}]_r| = H'(1) \leq \frac{n(D-1)D^n}{2}.$$

□

Corollary 5.6.5 below follows immediately from Lemmas 5.6.3 and 5.6.4.

**Corollary 5.6.5.** *Assume that  $\mathbf{f}$  is a regular sequence that satisfies Assumption (5.D). Then the degree of any minor of  $\mathcal{H}$  is bounded by  $n(D-1)D^n$ .*

**Example 5.6.6.** We consider again the system  $\mathbf{f} = (x_1^2 + x_2^2 - y_1, x_1x_2 + y_2x_2 + y_3x_1)$  in Example 5.4.3. Note that  $\mathbf{f}$  forms a regular sequence.

The Gröbner basis  $\mathcal{G}$  of  $\mathbf{f}$  with respect to the ordering  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$  is

$$\mathcal{G} = \{x_2^3 + y_3x_2^2 + (y_2^2 - y_1)x_2 + y_2y_3x_1 - y_1y_3, x_1^2 + x_2^2 - y_1, x_1x_2 + x_1y_3 + x_2y_2\}.$$

So,  $\mathbf{f}$  satisfies Assumption (5.D). The matrix with respect to the basis  $B_1 = \{1, x_2, x_1, x_2^2\}$  has the following degree pattern:

$$\begin{bmatrix} 0 & 1 & 1 & 2 \\ 1 & 2 & 2 & 3 \\ 1 & 2 & 2 & 3 \\ 2 & 3 & 3 & 4 \end{bmatrix}.$$

This degree pattern agrees with the result of Proposition 5.6.2. The determinant of this matrix is of degree 7, which is indeed smaller than  $n(D-1)D^n = 8$

Whereas, using the basis  $B_2 = \{1, x_2, x_2^2, x_2^3\}$  leads to another parametric Hermite matrix of different degrees. For  $1 \leq i, j \leq 4$ , the degree of its  $(i, j)$ -entry, which is equals to  $\text{trace}(\mathcal{L}_{x_2^{i+j-2}})$ , is bounded by  $\deg(x_2^{i-1}) + \deg(x_2^{j-1}) = i + j - 2$  using Proposition 5.6.2. Applying Lemma 5.6.3, the determinant is bounded by  $2 \sum_{i=0}^3 \deg(x_2^i) = 12$ .

By computing the parametric Hermite matrix of  $\mathbf{f}$  with respect to  $B_2$ , we obtain the degree pattern

$$\begin{bmatrix} 0 & 1 & 2 & 3 \\ 1 & 2 & 3 & 4 \\ 2 & 3 & 4 & 5 \\ 3 & 4 & 5 & 6 \end{bmatrix}$$

on its entries and a determinant of degree 11. Again, both of our theoretical bounds hold for this matrix.

**Remark 5.6.7.** Note that Assumption (5.D) requires a condition on the degrees of polynomials in the Gröbner basis  $\mathcal{G}$  of  $\langle \mathbf{f} \rangle$ . We remark that it is possible to establish similar bounds for the degrees of entries of our parametric Hermite matrix and its minors when the system  $\mathbf{f}$  satisfies a weaker property than Assumption (5.D) (we still keep the regularity assumption).

Indeed, we only need to assume that, for any  $g \in \mathcal{G}$ , the homogeneous component of the highest degree in  $\mathbf{x}$  of  $g$  does not depend on the parameters  $\mathbf{y}$ . Let  $D_{\mathbf{y}}$  be an upper bound of the partial degrees in  $\mathbf{y}$  of elements of  $\mathcal{G}$ . Under the change of variables  $x_i \mapsto x_i^{D_{\mathbf{y}}}$ ,  $\mathbf{f}$  is mapped to a new polynomial sequence that satisfies Assumption (5.D). Therefore, we easily deduce the two following bounds, which are similar to the ones of Proposition 5.6.2 and Corollary 5.6.5.

- $\deg(h_{i,j}) \leq D_{\mathbf{y}}(\deg(b_i) + \deg(b_j))$ ;
- The degree of any minor of  $\mathcal{H}$  is bounded by  $D_{\mathbf{y}} n(D-1)D^n$ .

Even though these bounds are not sharp anymore, they still allow us to compute the parametric Hermite matrices using evaluation & interpolation scheme and control the complexity of this computation in the instances for which Assumption (5.D) does not hold.

## 5.6.2 Complexity analysis of our algorithms

In this subsection, we analyze the complexity of our algorithms on generic systems.

Let  $\mathbf{f} = (f_1, \dots, f_n) \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  be a regular sequence, where  $\mathbf{y} = (y_1, \dots, y_t)$  and  $\mathbf{x} = (x_1, \dots, x_n)$ , satisfying Assumptions (5.A) and (5.D). We denote by  $\mathcal{G}$  the reduced Gröbner basis of  $\mathbf{f}$  with respect to the ordering  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$ . The basis  $\mathcal{B}$  is taken as all the monomials in  $\mathbf{x}$  that are irreducible by  $\mathcal{G}$ . Then,  $\mathcal{H}$  is the parametric Hermite matrix associated of  $\mathbf{f}$  with respect to  $\mathcal{B}$ .

We start by estimating the arithmetic complexity for computing the parametric Hermite matrix  $\mathcal{H}$  and its minors. We denote  $\lambda := n(D - 1)$  and  $\mathfrak{D} := n(D - 1)D^n$ .

**Proposition 5.6.8.** *Assume that  $\mathbf{f} = (f_1, \dots, f_n) \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  is a regular sequence that satisfies Assumptions (5.A) and (5.D). Let  $\delta$  be the dimension of the  $\mathbb{K}$ -vector space  $\mathbb{K}[\mathbf{x}]/\langle \mathbf{f} \rangle_{\mathbb{K}}$  where  $\mathbb{K} = \mathbb{Q}(\mathbf{y})$ . Let  $\mathcal{H}$  be the parametric Hermite matrix associated to  $\mathbf{f}$  constructed using  $\text{grevlex}(\mathbf{x})$  ordering. Then, by Lemma 5.4.5, the entries of the parametric Hermite matrix  $\mathcal{H}$  lie in  $\mathbb{Q}[\mathbf{y}]$ .*

Using the evaluation & interpolation scheme, one can compute  $\mathcal{H}$  within

$$O\left(\binom{t+2\lambda}{t} \left( n \binom{D+n+t}{n+t} + n^{\omega+2} D^{\omega n+1} + D^{(\omega+1)n} \right)\right)$$

arithmetic operations in  $\mathbb{Q}$ , where, by Bézout's bound,  $\delta$  is bounded by  $D^n$ .

Moreover, each minor (including the determinant) of  $\mathcal{H}$  can be computed using

$$O\left(\binom{t+\mathfrak{D}}{t} \left( \delta^2 \binom{t+2\lambda}{t} + \delta^\omega \right)\right)$$

arithmetic operations in  $\mathbb{Q}$ .

*Proof.* By Lemma 5.6.4 and Proposition 5.6.2, the highest degree among the entries of the Hermite matrix  $\mathcal{H}$  is bounded by  $2\lambda = 2n(D - 1)$ . The evaluation & interpolation scheme of Subsection 5.4.4 requires computing  $\binom{t+2\lambda}{t}$  specialized Hermite matrices. We first analyze the complexity for computing each of those specialized Hermite matrices.

The evaluation of  $\mathbf{f}$  at each point  $\eta \in \mathbb{Q}^t$  costs  $O\left(n \binom{D+n+t}{n+t}\right)$  arithmetic operations in  $\mathbb{Q}$ .

As the highest degree in the Gröbner basis of  $\mathbf{f}(\eta, \cdot)$  with respect to the  $\text{grevlex}(\mathbf{x})$  ordering is bounded by  $n(D - 1)$ , the computation of this Gröbner basis can be done within  $O(nD^{\omega n})$  arithmetic operations in  $\mathbb{Q}$  (see [58, Theorem 5.1]).

Next, we compute the matrices representing the  $\mathcal{L}_{x_i}$ 's. Using [58, Algo. 4], we obtain an arithmetic complexity of  $O(Dn^{\omega+2}\delta^\omega)$  ([58, Theorem 5.1]) for computing such  $n$  matrices, where  $\omega$  is the exponential constant for matrix multiplication.

The traces of these matrices are then computed using  $n\delta$  additions in  $\mathbb{Q}$ . The subroutine `BMatrices` consists of essentially  $\delta$  multiplication of  $\delta \times \delta$  matrices (with entries in  $\mathbb{Q}$ ). This leads to an arithmetic complexity  $O(\delta^{\omega+1})$ . Next, the computation of each entry  $h_{i,j}$  is simply a vector multiplication of length  $\delta$ , whose complexity is  $O(\delta)$ . Doing so for  $\delta^2$  entries, `TraceComputing` takes in overall  $O(\delta^3)$  arithmetic operations in  $\mathbb{Q}$ .

Thus, as  $\delta \leq D^n$ , the complexity of the evaluation step lies in

$$O\left(\binom{t+2\lambda}{t} \left( n \binom{D+n+t}{n+t} + n^{\omega+2} D^{\omega n+1} + D^{(\omega+1)n} \right)\right).$$

Finally, we interpolate  $\delta^2$  entries which are polynomials in  $\mathbb{Q}[\mathbf{y}]$  of degree at most  $2\lambda$ . Using the multivariate interpolation algorithm of [36], the complexity of this step therefore lies in  $O\left(\delta^2 \binom{t+2\lambda}{t} \log^2 \binom{t+2\lambda}{t} \log \log \binom{t+2\lambda}{t}\right)$ .

Summing up the both steps, we conclude that the parametric Hermite matrix  $\mathcal{H}$  can be obtained within

$$O\sim\left(\binom{t+2\lambda}{t} \left( n \binom{D+n+t}{n+t} + n^{\omega+2} D^{\omega n+1} + D^{(\omega+1)n} \right)\right)$$

arithmetic operations in  $\mathbb{Q}$ .

Similarly, the minors of  $\mathcal{H}$  can be computed using the technique of evaluation & interpolation. By Corollary 5.6.5, the degree of every minor of  $\mathcal{H}$  is bounded by  $\mathfrak{D}$ . We specialize  $\mathcal{H}$  at  $\binom{t+\mathfrak{D}}{t}$  points in  $\mathbb{Q}^t$  and compute the corresponding minor of each specialized Hermite matrix. This step takes

$$O\left(\binom{t+\mathfrak{D}}{t} \left( \delta^2 \binom{t+2\lambda}{t} + \delta^\omega \right)\right)$$

arithmetic operations in  $\mathbb{Q}$ . Finally, using the multivariate interpolation algorithm of [36], it requires

$$O\left(\binom{t+\mathfrak{D}}{t} \log^2 \binom{t+\mathfrak{D}}{t} \log \log \binom{t+\mathfrak{D}}{t}\right)$$

arithmetic operations in  $\mathbb{Q}$  to interpolate the final minor. Therefore, the whole complexity for computing each minor of  $\mathcal{H}$  lies within

$$O\sim\left(\binom{t+\mathfrak{D}}{t} \left( \delta^2 \binom{t+2\lambda}{t} + \delta^\omega \right)\right).$$

□

Finally, we state our main result, which is Theorem 5.1.3 below. It estimates the arithmetic complexity of Algorithms 5.3 and 5.4.

**Theorem 5.1.3.** *Let  $\mathbf{f} \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  be a regular sequence such that the ideal  $\langle \mathbf{f} \rangle$  is radical and  $\mathbf{f}$  satisfies Assumptions (5.A) and (5.D). Recall that  $\mathfrak{D}$  denotes  $n(D-1)D^n$ . Then, we have the following statements:*

i) The arithmetic complexity of Algorithm 5.3 lies in

$$O\sim\left(\binom{t+\mathfrak{D}}{t} 2^{3t} n^{2t+1} D^{2nt+n+2t+1}\right).$$

ii) Algorithm 5.4, which is probabilistic, computes a set of semi-algebraic descriptions solving Problem (RRC) within

$$O\sim\left(\binom{t+\mathfrak{D}}{t} 2^{3t} n^{2t+1} D^{3nt+2(n+t)+1}\right)$$

arithmetic operations in  $\mathbb{Q}$  in case of success.

iii) The semi-algebraic descriptions output by Algorithm 5.4 consist of polynomials in  $\mathbb{Q}[\mathbf{y}]$  of degree bounded by  $\mathfrak{D}$ .

*Proof.* As Assumption (5.D) holds, we have that  $w_\infty = 1$  and  $w_{\mathcal{H}}$  is the square-free part of  $\det(\mathcal{H})$ .

Therefore, after computing the parametric Hermite matrix  $\mathcal{H}$  and its determinant, whose complexity is given by Proposition 5.6.8, Algorithm 5.3 essentially consists of computing sample points of the connected components of the algebraic set  $\mathbb{R}^t \setminus V(\det(\mathcal{H}))$ .

By Corollary 5.6.5, the degree of  $\det(\mathcal{H})$  is bounded by  $\mathfrak{D}$ . Applying Corollary 5.2.3, we obtain the following arithmetic complexity for this computation of sample points

$$O\sim\left(\binom{t+\mathfrak{D}}{t} 2^{3t} \mathfrak{D}^{2t+1}\right) \simeq O\sim\left(\binom{t+\mathfrak{D}}{t} 2^{3t} n^{2t+1} D^{2nt+n+2t+1}\right).$$

Also by Corollary 5.2.3, the finite subset of  $\mathbb{Q}^t$  output by SamplePoints has cardinal bounded by  $2^t \mathfrak{D}^t$ . Thus, evaluating the specializations of  $\mathcal{H}$  at those points and their signatures costs in total  $O\left(2^t \mathfrak{D}^t \left(\delta^2 \binom{2\lambda+t}{t} + \delta^{\omega+1/2}\right)\right)$  arithmetic operations in  $\mathbb{Q}$  using [9, Algorithm 8.43].

Therefore, the complexity of SamplePoints dominates the whole complexity of the algorithm. We conclude that Algorithm 5.3 runs within

$$O\sim\left(\binom{t+\mathfrak{D}}{t} 2^{3t} n^{2t+1} D^{2nt+n+2t+1}\right)$$

arithmetic operations in  $\mathbb{Q}$ .

For Algorithm 5.4, we start by choosing randomly a matrix  $Q$  and compute the matrix  $\mathcal{H}_Q = Q^T \cdot \mathcal{H} \cdot Q$ . Then, we compute the leading principal minors  $M_1, \dots, M_\delta$  of  $\mathcal{H}_Q$ . Using Proposition 5.6.8, this step admits the arithmetic complexity bound

$$O\sim\left(\delta \binom{t+\mathfrak{D}}{t} \left(\delta^2 \binom{t+2\lambda}{t} + \delta^\omega + \log^2 \binom{t+\mathfrak{D}}{t}\right)\right).$$

Next, Algorithm 5.4 computes sample points for the connected components of the semi-algebraic set defined by  $\bigwedge_{i=1}^{\delta} M_i \neq 0$ . Since the degree of each  $M_i$  is bounded by  $\mathfrak{D}$ , Corollary 5.2.3 gives the arithmetic complexity

$$O\left(\binom{t + \mathfrak{D}}{t} D^{nt+n} 2^{3t} \mathfrak{D}^{2t+1}\right) \simeq O\left(\binom{t + \mathfrak{D}}{t} 2^{3t} n^{2t+1} D^{3nt+2(n+t)+1}\right).$$

It returns a finite subset of  $\mathbb{Q}^t$  whose cardinal is bounded by  $(2\delta\mathfrak{D})^t$ . The evaluation of the leading principal minors' sign patterns at those points has the arithmetic complexity lying in  $O(2^t \delta^{t+1} \mathfrak{D}^{2t}) \simeq O(2^t n^{2t} D^{3nt+n+2t})$ .

Again, the complexity of SamplePoints dominates the whole complexity of Algorithm 5.4. The proof of Theorem 5.1.3 is then finished.  $\square$

**Probability aspect.** Here, we give some short remarks on the probabilistic aspect of Algorithms 5.3 and 5.4.

These two algorithms rely on the geometric resolution algorithm of [87] for solving zero-dimensional systems appearing in the computation of sample points per connected components described in Section 5.2. Recall that the geometric resolution is a probabilistic algorithm, which makes various random choices (changes of variables, generic points to specialize) to ensure certain properties of the intermediate systems.

As explained in [87], the bad choices are enclosed in strict algebraic subsets of certain affine spaces, which implies that almost any set of random choices leads to a correct result. In general, even though one can check whether the points output by geometric resolution are solutions of the input system, some solutions can be missing. Thus, the geometric resolution is not Las Vegas. It is worth note that, by replacing the geometric resolution algorithm by an algorithm for solving zero-dimensional system using Gröbner basis, we obtain a deterministic version of the subroutine SamplePoints.

Besides, Algorithm 5.4 depends also on the choice of the matrix  $Q$ . By Lemma 5.5.8, any choice of  $Q$  from a prescribed dense Zariski open subset of  $\text{GL}(n, \mathbb{C})$  will work. As the purpose of choosing  $Q$  is to ensure that none of the leading principal minors of  $Q^T \cdot \mathcal{H} \cdot Q$  are identically zero. One can check easily whether a good matrix  $Q$  is found. Again, this can be made deterministic.

## 5.7 Practical implementation & Experimental results

### 5.7.1 Remark on the implementation of Algorithm 5.4

Recall that Algorithm 5.4 leads us to compute sample points per connected components of the non-vanishing set of the leading principal minors  $(M_1, \dots, M_\delta)$ . Comparing to Algorithm 5.3 in which we only compute sample points for  $\mathbb{R}^t \setminus V(M_\delta)$ , the complexity of Algorithm 5.4 contains an extra factor of  $D^{nt}$  due to the higher number of polynomials given as input to the subroutine SamplePoints. Even though the complexity bounds of these two algorithms both lie in  $D^{O(nt)}$ , the

extra factor  $D^{nt}$  mentioned above sometimes becomes the bottleneck of Algorithm 5.4 for tackling practical problems. Therefore, we introduce the following optimization in our implementation of Algorithm 5.4.

We start by following exactly the steps (1-4) of Algorithm 5.4 to obtain the leading principal minors  $(M_1, \dots, M_\delta)$  and the polynomial  $w_\infty$ . Then, by calling the subroutine SamplePoints on the input  $M_\delta \neq 0 \wedge w_\infty \neq 0$ , we compute a set of sample points (and their corresponding numbers of real roots)  $\{(\eta_1, r_1), \dots, (\eta_\ell, r_\ell)\}$  that solves the weak-version of Problem (RRC). We obtain from this output all the possible numbers of real roots that the input system can admit.

For each value  $0 \leq r \leq \delta$ , we define

$$\Phi_r = \{\sigma = (\sigma_1, \dots, \sigma_\delta) \in \{-1, 1\}^\delta \mid \text{the sign variation of } \sigma \text{ is } (\delta - r)/2\}.$$

If  $r \not\equiv \delta \pmod{2}$ ,  $\Phi_r = \emptyset$ .

For  $\sigma \in \Phi_r$  and  $\eta \in \mathbb{R}^t \setminus V(w_\infty)$  such that  $\text{sign}(M_i(\eta)) = \sigma_i$  for every  $1 \leq i \leq \delta$ , the signature of  $\mathcal{H}(\eta)$  is  $r$ . As a consequence, for any  $\eta$  in the semi-algebraic set defined by

$$(w_\infty \neq 0) \wedge (\bigvee_{\sigma \in \Phi_r} (\bigwedge_{i=1}^\delta \text{sign}(M_i) = \sigma_i)),$$

the system  $\mathbf{f}(\eta, \cdot)$  has exactly  $r$  distinct real solutions.

Therefore,  $(\mathcal{S}_{r_i})_{1 \leq i \leq \ell}$  is a collection of semi-algebraic sets solving Problem (RRC). Then, we can simply return  $\{(\Phi_{r_i}, \eta_i, r_i) \mid 1 \leq i \leq \ell\}$  as the output of Algorithm 5.4 without any further computation. Note that, by doing so, we may return sign conditions which are not realizable.

We discuss now about the complexity aspect of the steps described above. For  $r \equiv \delta \pmod{2}$ , the cardinal of  $\Phi_r$  is  $\binom{\delta}{(\delta-r-2)/2}$ . In theory, the total cardinal of all the  $\Phi_{r_i}$ 's ( $1 \leq i \leq \ell$ ) can go up to  $2^{\delta-1}$ , which is doubly exponential in the number of variables  $n$ . However, in the instances that are actually tractable by the current state of the art,  $2^\delta$  is still smaller than  $\delta^{3t}$ . And when it is the case, following this approach has better performance than computing the sample points of the semi-algebraic set defined by  $\bigwedge_{i=1}^\delta M_i \neq 0$ . Otherwise, when  $2^\delta$  exceeds  $\delta^{3t}$ , we switch back to the computation of sample points.

This implementation of Algorithm 5.4 does not change the complexity bound given in Theorem 5.1.3.

## 5.7.2 Experiments

This subsection reports on the practical performance of several real root classification algorithms on various test instances and applications.

The computation is carried out on a computer of Intel(R) Xeon(R) CPU E7-4820 2GHz and 1.5 TB of RAM. The timings are given in seconds (s.), minutes (m.) and hours (h.). The symbol  $\infty$  means that the computation cannot finish within 240 hours.

We implement Algorithm 5.3 and 5.4 in Maple. This implementation relies on the library FGB [66] for carrying out the Gröbner basis computation required to compute Hermite matrices (Algorithm 5.2). It also calls to the library RAGLIB to compute sample points of semi-algebraic

sets, which makes use of the library `MSOLVE` for solving zero-dimensional systems appearing in those computations.

Throughout this subsection, the column `HERMITE` reports on the computational data of our algorithms based on parametric Hermite matrices described in Section 5.5. It uses the notations below:

- `MAT`: the timing for computing a parametric Hermite matrix  $\mathcal{H}$ .
- `DET`: the runtime for computing the determinant of  $\mathcal{H}$ .
- `MIN`: the timing for computing the leading principal minors of  $\mathcal{H}$ .
- `SP`: the runtime for computing at least one points per each connected component of the semi-algebraic set  $\mathbb{R}^t \setminus V(\det(\mathcal{H}))$ .
- `DEG`: the highest degree among the leading principal minors of  $\mathcal{H}$ .

**Generic systems.** In this paragraph, we report on the results obtained with generic inputs, i.e., randomly chosen dense polynomials  $(f_1, \dots, f_n) \subset \mathbb{Q}[y_1, \dots, y_t][x_1, \dots, x_n]$ . The total degrees of input polynomials are given as a list  $D = [\deg(f_1), \dots, \deg(f_n)]$ .

We first compare the algorithms using Hermite matrices (Section 5.5) with the Sturm-based algorithm (Section 5.3) for solving Problem (RRC). The column `STURM` of Fig. (5.2) shows the experimental results of the Sturm-based algorithm. It contains the following sub-columns:

- `ELIM`: the timing for computing the eliminating polynomial.
- `SRES`: the timing for computing the subresultant coefficients in the Sturm-based algorithm.
- `SP-S`: the timing for computing points per connected component of the non-vanishing set of the last subresultant coefficient.
- `DEG-S`: the highest degree among the subresultant coefficients.

We observe that the sum of `MAT-H` and `MIN-H` is smaller than the sum of `ELIM` and `SRES`. Hence, obtaining the input for the sample point computation in `HERMITE` strategy is easier than in `STURM` strategy. We also remark that the degree `DEG-H` is much smaller than `DEG-S`, that explains why the computation of sample points using Hermite matrices is faster than using the subresultant coefficients.

We conclude that the parametric Hermite matrix approach outperforms the Sturm-based one both on the timings and the degree of polynomials in the output formulas.

In Fig. (5.3), we compare our algorithms using parametric Hermite matrices with two Maple packages for solving parametric polynomial systems: `ROOTFINDING[PARAMETRIC]` [76] and `REGULARCHAINS[PARAMETRICSYSTEMTOOLS]` [202]. The new notations used in Fig. (5.3) are explained below.

$t$	$D$	HERMITE					STURM				
		MAT	MIN	SP	total	DEG	ELIM	SRES	SP-S	total	DEG-S
2	[2, 2]	.07 s	.01 s	.3 s	.4 s	8	.01 s	.1 s	2 s	2.2 s	12
2	[3, 2]	.1 s	.12 s	4.8 s	5 s	18	.05 s	.5 s	15 s	16 s	30
2	[2, 2, 2]	.3 s	.3 s	33 s	34 s	24	.08 s	2 s	8 m	8 m	56
2	[3, 3]	.3 s	.8 s	3 m	3 m	36	.1 s	3 s	20 m	20 m	72
3	[2, 2]	.1 s	.02 s	26 s	27 s	8	.07 s	.1 s	40 s	40 s	12
3	[3, 2]	.2 s	.2 s	3 h	3 h	18	.1 s	1 s	$\infty$	$\infty$	30
3	[2, 2, 2]	.5 s	7 s	32 h	32 h	24	.15 s	10 m	$\infty$	$\infty$	56
3	[4, 2]	.6 s	12 s	90 h	90 h	32	.2 s	12 m	$\infty$	$\infty$	56
3	[3, 3]	1 s	27 s	$\infty$	$\infty$	36	.2 s	15 m	$\infty$	$\infty$	72

Figure 5.2: Generic random dense systems

- The column RF stands for the ROOTFINDING[PARAMETRIC] package. To solve a parametric polynomial systems, it consists of computing a discriminant variety  $\mathcal{D}$  and then computing an open CAD of  $\mathbb{R}^t \setminus \mathcal{D}$ . *This package does not return explicit semi-algebraic formulas but an encoding based on the real roots of some polynomials.*

This column contains:

- DV : the runtime of the command DISCRIMINANTVARIETY that computes a set of polynomials defining a discriminant variety  $\mathcal{D}$  associated to the input system.
- CAD : the runtime of the command CELLEDECOMPOSITION that outputs semi-algebraic formulas by computing an open CAD for the semi-algebraic set  $\mathbb{R}^t \setminus \mathcal{D}$ .
- The column RC stands for Maple’s library REGULARCHAINS[PARAMETRICSYSTEMTOOLS]. The algorithms implemented in this package is given in [202]. It also contains two sub-columns:
  - BP : the runtime of the command BORDERPOLYNOMIAL that returns a set of polynomials.
  - RRC : the runtime of the command REALROOTCLASSIFICATION. We call this command with the option `output='samples'` to compute at least one point per connected component of the complementary of the real algebraic set defined by border polynomials.

Note that, in a strategy for solving the weak-version of Problem (RRC), DISCRIMINANTVARIETY and BORDERPOLYNOMIAL can be completely replaced by parametric Hermite matrices.

On generic systems, the determinant of our parametric Hermite matrix coincides with the output of DISCRIMINANTVARIETY, which we denote by  $w$ . Whereas, because of the elimination BORDERPOLYNOMIAL returns several polynomials, one of them is  $w$ .

In Fig. (5.3), the timings for computing a parametric Hermite matrix is negligible. Comparing the columns DET, DV and BP, we remark that the time taken to obtain  $w$  through the determinant

of parametric Hermite matrices is much smaller than using DISCRIMINANTVARIETY or BORDERPOLYNOMIAL.

For computing the polynomial  $w$ , using parametric Hermite matrices allows us to reach the instances that are out of reach of DISCRIMINANTVARIETY, for example, the instances  $\{t = 3, D = [2, 2, 2]\}$ ,  $\{t = 3, D = [4, 2]\}$ ,  $\{t = 3, D = [3, 3]\}$  and  $\{t = 4, D = [2, 2]\}$  in Fig. (5.3) below. Moreover, we succeed to compute the semi-algebraic formulas for  $\{t = 3, D = [2, 2, 2]\}$ ,  $\{t = 3, D = [4, 2]\}$  and  $\{t = 4, D = [2, 2]\}$ . Using the implementation in Subsection 5.7.1, we obtain the semi-algebraic formulas of degrees bounded by  $\deg(w)$ .

Therefore, for these generic systems, our algorithm based on parametric Hermite matrices outperforms DISCRIMINANTVARIETY and BORDERPOLYNOMIAL for obtaining a polynomial that defines the boundary of semi-algebraic sets over which the number of real solutions are invariant. Moreover, using the minors of parametric Hermite matrices, we can compute semi-algebraic formulas of problems that are out of reach of CELLECOMPOSITION and REALROOTCLASSIFICATION.

$t$	$d$	HERMITE					DEG	RF			RC		
		MAT	DET	SP	total	DV		CAD	total	BP	RRC	total	
2	[2, 2]	.07 s	.01 s	.3 s	.4 s	8	.1 s	.3 s	.4 s	.1 s	1 s	1.1 s	
2	[3, 2]	.1 s	.2 s	4.8 s	5 s	18	1 m	5 s	1 m	.3 s	12 s	12 s	
2	[2, 2, 2]	.3 s	.3 s	33 s	34 s	24	17m	32 s	17m	23 s	2 m	2 m	
2	[3, 3]	.3 s	.8 s	3 m	3 m	36	2 h	4 m	2 h	8 s	4 m	4 m	
3	[2, 2]	.1 s	.02 s	26 s	27 s	8	1 s	35 s	36 s	.2 s	12m	12m	
3	[3, 2]	.2 s	.2 s	3 h	3 h	18	2 h	84 h	86 h	3 s	37 h	37 h	
3	[2, 2, 2]	.5 s	7 s	32 h	32 h	24	$\infty$	$\infty$	$\infty$	20m	$\infty$	$\infty$	
3	[4, 2]	.6 s	12 s	90 h	90 h	32	$\infty$	$\infty$	$\infty$	12m	$\infty$	$\infty$	
3	[3, 3]	.7 s	27 s	$\infty$	$\infty$	36	$\infty$	$\infty$	$\infty$	15m	$\infty$	$\infty$	
4	[2, 2]	.2 s	.1 s	8 m	8 m	8	4 s	$\infty$	$\infty$	1 s	$\infty$	$\infty$	

Figure 5.3: Generic random dense systems

Especially, since the polynomials in our outputs are obtained as minors of parametric Hermite matrices, these matrices provide a compact determinantal representation of the output formulas, which then facilitates their evaluation. We illustrate this claim by reporting in Table 5.4 on the timings of these two different tasks for 1000 points  $\eta$ :

- Evaluating the signature of  $\mathcal{H}(\eta)$  (the column SIGN);
- Evaluating the principal minors of  $\mathcal{H}$  (the column MINORS);
- Solving specialized systems  $\mathbf{f}(\eta, \cdot)$  using MSOLVE, FGB and ROOTFINDING[ISOLATE] of Maple (the columns MSOLVE, FGB and ISOLATE).

We note that evaluating the signatures of specialized Hermite matrices is faster than evaluating the minors. On the other hand, solving a specialized system would depend strongly on the

System	$t$	$n$	$D$	SIGN	MINORS	MSOLVE	FGB	ISOLATE
Dense	2	2	[2, 2]	.5 s	.2 s	2 s	12 s	33 s
	2	3	[2, 2, 2]	2 s	4 s	5 s	15 s	110 s
	2	2	[3, 3]	3 s	6 s	4 s	12 s	65 s
	2	2	[5, 2]	7 s	18 s	5 s	14 s	55 s
	2	2	[4, 3]	10 s	30 s	6 s	15 s	80 s
Dense	2	2	[2, 2]	.8 s	.4 s	2 s	10 s	16 s
	2	3	[2, 2, 2]	6 s	30 s	5 s	15 s	80 s
	2	2	[3, 3]	9 s	90 s	4 s	12 s	65 s

Figure 5.4: Timings for evaluating the formulas.

number of variables  $n$  while evaluating the signatures depends on the number of parameters  $t$ . In the above examples where  $n = 2$  and  $t = 3$ , solving the specialized systems is better. Even though, the only library for solving polynomial systems is faster than evaluating the signatures on these examples is MSOLVE, which is highly optimized in C.

In what follows, we consider the systems coming from some applications as test instances. These examples allow us to observe the behavior of our algorithms on non-generic systems.

**Kuramoto model.** This application is introduced in [126], which is a dynamical system used to model synchronization among some given coupled oscillators. Here we consider only the model constituted by 4 oscillators. The maximum number of real solutions of steady-state equations of this model was an open problem before it is solved in [99] using numerical homotopy continuation methods. However, to the best of our knowledge, there is no exact algorithm that is able to solve this problem. We present in what follows the first solution using symbolic computation. Moreover, our algorithm can return the semi-algebraic formulas defining the regions over which the number of real solutions is invariant.

As explained in [99], we consider the system  $\mathbf{f}$  of the following equations

$$\begin{cases} y_i - \sum_{j=1}^4 (s_i c_j - s_j c_i) = 0 \\ s_i^2 + c_i^2 = 1 \end{cases} \text{ for } 1 \leq i \leq 3,$$

where  $(s_1, s_2, s_3)$  and  $(c_1, c_2, c_3)$  are variables and  $(y_1, y_2, y_3)$  are parameters. We are asked to compute the maximum number of real solutions of  $\mathbf{f}(\eta, \cdot)$  when  $\eta$  varies over  $\mathbb{R}^3$ . This leads us to solve the weak version of Problem (RRC) for this parametric system.

We first construct the parametric Hermite matrix  $\mathcal{H}$  associated to this system. This matrix <sup>1</sup>

<sup>1</sup>The matrix is available at <https://github.com/huuphuocle/Kuramoto4>.

is of size  $14 \times 14$  and has the following degree pattern:

$$\begin{bmatrix} 0 & 3 & 3 & 0 & 0 & 0 & 6 & 4 & 3 & 3 & 6 & 3 & 4 & 9 \\ 3 & 6 & 6 & 3 & 3 & 3 & 9 & 7 & 6 & 6 & 9 & 6 & 7 & 12 \\ 3 & 6 & 6 & 3 & 3 & 3 & 9 & 7 & 6 & 6 & 9 & 6 & 7 & 12 \\ 0 & 3 & 3 & 2 & 2 & 2 & 6 & 4 & 5 & 5 & 6 & 5 & 4 & 9 \\ 0 & 3 & 3 & 2 & 2 & 2 & 6 & 4 & 5 & 5 & 6 & 5 & 4 & 9 \\ 0 & 3 & 3 & 2 & 2 & 2 & 6 & 4 & 5 & 5 & 6 & 5 & 4 & 9 \\ 6 & 9 & 9 & 6 & 6 & 6 & 12 & 10 & 9 & 9 & 12 & 9 & 10 & 15 \\ 4 & 7 & 7 & 4 & 4 & 4 & 10 & 8 & 7 & 7 & 10 & 7 & 8 & 13 \\ 3 & 6 & 6 & 5 & 5 & 5 & 9 & 7 & 8 & 8 & 9 & 8 & 7 & 12 \\ 3 & 6 & 6 & 5 & 5 & 5 & 9 & 7 & 8 & 8 & 9 & 8 & 7 & 12 \\ 6 & 9 & 9 & 6 & 6 & 6 & 12 & 10 & 9 & 9 & 12 & 9 & 10 & 15 \\ 3 & 6 & 6 & 5 & 5 & 5 & 9 & 7 & 8 & 8 & 9 & 8 & 7 & 12 \\ 4 & 7 & 7 & 4 & 4 & 4 & 10 & 8 & 7 & 7 & 10 & 7 & 8 & 13 \\ 9 & 12 & 12 & 9 & 9 & 9 & 15 & 13 & 12 & 12 & 15 & 12 & 13 & 18 \end{bmatrix}.$$

The polynomial  $w_\infty$  has the factors  $y_1 + y_2$ ,  $y_2 + y_3$ ,  $y_3 + y_1$  and  $y_1 + y_2 + y_3$ . The polynomial  $w_{\mathcal{H}}$  has degree 48 (c.f. [99]). We denote by  $w$  the polynomial  $w_\infty \cdot w_{\mathcal{H}}$ .

Note that the polynomial system has real roots only if  $|y_i| \leq 3$  (c.f. [99]). So we only need to consider the compact connected components of  $\mathbb{R}^3 \setminus V(w)$ . Since the polynomial  $w$  is invariant under any permutation acting on  $(y_1, y_2, y_3)$ , we exploit this symmetry to accelerate the computation of sample points.

Following the critical point method, we compute the critical points of the map  $(y_1, y_2, y_3) \mapsto y_1 + y_2 + y_3$  restricted to  $\mathbb{R}^3 \setminus V(w)$ ; this map is also symmetric. We apply the change of variables

$$(y_1, y_2, y_3) \mapsto (e_1, e_2, e_3),$$

where  $e_1 = y_1 + y_2 + y_3$ ,  $e_2 = y_1y_2 + y_2y_3 + y_3y_1$  and  $e_3 = y_1y_2y_3$  are elementary symmetric polynomials of  $(y_1, y_2, y_3)$ . This change of variables reduces the number of distinct solutions of zero-dimensional systems involved in the computation and, therefore, reduces the computation time.

From the sample points obtained by this computation, we derive the possible number of real solutions and conclude that the system  $f$  has at most 10 distinct real solutions when  $(y_1, y_2, y_3)$  varies over  $\mathbb{R}^3 \setminus V(w)$ . This agrees with the result given in [99]. We show below a list of parameter values such that the system has respectively 2, 4, 6, 8 and 10 distinct real solutions.

Number of solutions	$(y_1, y_2, y_3)$
2 solutions	$[-2, -0.03, 0.22]$
4 solutions	$[1, -0.09, 0.16]$
6 solutions	$[0, -0.7, -0.48]$
8 solutions	$[0.08, -0.03, 0.22]$
10 solutions	$\left[ \begin{array}{ccc} 274945023031 & -68723139707 & -549808278091 \\ 2199023255552 & 549755813888 & 4398046511104 \end{array} \right]$

Fig. (5.5) reports on the timings for computing the parametric Hermite matrix (MAT), for computing its determinant (DET) and for computing the sample points (SP). We stop both of the commands DISCRIMINANTVARIETY and BORDERPOLYNOMIAL after 240 hours without obtaining the polynomial  $w$ .

MAT	HERMITE		total	DV	BP
	DET	SP			
2 m	1 h	85 h	86 h	$\infty$	$\infty$

Figure 5.5: Kuramoto model for 4 oscillators

**Static output feedback.** The second non-generic example comes from the problem of static output feedback [105]. Given the matrices  $A \in \mathbb{R}^{\ell \times \ell}$ ,  $B \in \mathbb{R}^{\ell \times 2}$ ,  $C \in \mathbb{R}^{1 \times \ell}$  and a parameter vector  $P = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \in \mathbb{R}^2$ , the characteristic polynomial of  $A + BPC$  writes

$$f(s, \mathbf{y}) = \det(sI_\ell - A - BKC) = f_0(s) + y_1 f_1(s) + y_2 f_2(s),$$

where  $s$  is a complex variable.

We want to find a matrix  $P$  such that all the roots of  $f(s, \mathbf{y})$  must lie in the open left half-plane. By substituting  $s$  by  $x_1 + ix_2$ , we obtain the following system of real variables  $(x_1, x_2)$  and parameters  $(y_1, y_2)$ :

$$\begin{cases} \Re(f(x_1 + ix_2, \mathbf{y})) = 0 \\ \Im(f(x_1 + ix_2, \mathbf{y})) = 0 \\ x_1 < 0 \end{cases}$$

Note that the total degree of these equations equals  $\ell$ .

We are now interested in solving the weak-version of Problem (RRC) on the system  $\Re(f) = \Im(f) = 0$ . We observe that this system satisfies Assumptions (5.A) and (5.C). Let  $\mathcal{H}$  be the parametric Hermite matrix  $\mathcal{H}$  of this system with respect to the usual basis we consider in this paper. This matrix  $\mathcal{H}$  behaves very differently from generic systems.

Computing the determinant of  $\mathcal{H}$  (which is an element of  $\mathbb{Q}[\mathbf{y}]$ ) and taking its square-free part allows us to obtain the same output  $w$  as DISCRIMINANTVARIETY. However, this direct approach appears to be very inefficient as the determinant appears as a large power of the output polynomial.

For example, for a value  $\ell$ , we observe that the system consists of two polynomials of degree  $\ell$ . The determinant of  $\mathcal{H}$  appears as  $w^{2\ell}$ , where  $w$  has degree  $2(\ell - 1)$ . The bound we establish on the degree of this determinant is  $2(\ell - 1)\ell^2$ , which is much larger than what happens in this case. Therefore, we need to introduce the optimization below to adapt our implementation of Algorithm 5.3 to this problem.

We observe that, on these examples, the polynomial  $w$  can be extracted from a smaller minor instead of computing the determinant  $\mathcal{H}$ . To identify such a minor, we reduce  $\mathcal{H}$  to a matrix whose entries are univariate polynomials with coefficients lying in a finite field  $\mathbb{Z}/p\mathbb{Z}$  as follow.

Let  $u$  be a new variable. We substitute each  $y_i$  by random linear forms in  $\mathbb{Q}[u]$  in  $\mathcal{H}$  and then compute  $\mathcal{H} \bmod p$ . Then, the matrix  $\mathcal{H}$  is turned into a matrix  $\mathcal{H}_u$  whose entries are elements of  $\mathbb{Z}/p\mathbb{Z}[u]$ . The computation of the leading principal minors of  $\mathcal{H}_u$  is much easier than the one of  $\mathcal{H}$  since it involves only univariate polynomials and does not suffer from the growth of bit-sizes as for the rational numbers.

Next, we compute the sequence of the leading principal minors of  $\mathcal{H}_u$  in decreasing order, starting from the determinant. Once we obtain a minor, of some size  $r$ , that is not divisible by  $\overline{w}_u$ , we stop and take the index  $r+1$ . Then, we compute the square-free part of the  $(r+1) \times (r+1)$  leading principal minor of  $\mathcal{H}$ , which can be done through evaluation-interpolation method. This yields a Monte Carlo implementation that depends on the choice of the random linear forms in  $\mathbb{Q}[u]$  and the finite field to compute the polynomial  $w$ .

In Fig. (5.6), we report on some computational data for the static output feedback problem. Here we choose the prime  $p$  to be 65521 so that the elements of the finite field  $\mathbb{Z}/p\mathbb{Z}$  can be represented by a machine word of 32 bits. We consider different values of  $\ell$  and the matrices  $A, B, C$  are chosen randomly. On these examples, our algorithm returns the same output as the one of DISCRIMINANTVARIETY. Whereas, BORDERPOLYNOMIAL (BP) returns a list of polynomials which contains our output and other polynomials of higher degree.

The timings of our algorithm are given by the two following columns:

- The column MAT shows the timings for computing parametric Hermite matrices  $\mathcal{H}$ .
- The column COMP-W shows the timings for computing the polynomials  $w$  from  $\mathcal{H}$  using the strategy described as above.

We observe that our algorithm (MAT + COMP-W) wins some constant factor comparing to DISCRIMINANTVARIETY (DV). On the other hand, BORDERPOLYNOMIAL (BP) performs less efficiently than the other two algorithms in these examples.

Since the degrees of the polynomials  $w$  here (given as DEG-W) are small comparing with the bounds in the generic case. Hence, unlike the generic cases, the computation of the sample points in these problems is negligible as being reported in the column SP.

$\ell$	MAT	HERMITE COMP-W	total	DV	BP	SP	DEG-W
	5	2 s	1 s	3 s	30 s	1.5 m	.2 s
6	12 s	5 s	17 s	90 s	30 m	.4 s	10
7	1 m	6 m	7 m	16 m	4 h	1 s	12
8	4 m	50 m	1 h	1.5 h	34 h	3 s	14

Figure 5.6: Static output feedback

## Chapter 6

# One block quantifier elimination for regular polynomial systems of equations

**Abstract.** Quantifier elimination over the reals is a central problem in computational real algebraic geometry, polynomial system solving and symbolic computation. Given a semi-algebraic formula (whose atoms are polynomial constraints) with quantifiers on some variables, it consists in computing a logically equivalent formula involving only unquantified variables. When there is no alternation of quantifiers, one has a *one block* quantifier elimination problem.

We study in this chapter a variant of the one block quantifier elimination in which we compute an almost equivalent formula of the input.

Our main contribution is a new probabilistic efficient algorithm for solving this variant when the input is a system of polynomial equations satisfying some regularity assumptions. When the input is generic, involves  $s$  polynomials of degree bounded by  $D$  with  $n$  quantified variables and  $t$  unquantified ones, we prove that this algorithm outputs semi-algebraic formulas of degree bounded by  $\mathfrak{B}$  using

$$O\left(8^t \mathfrak{B}^{3t+2} \binom{t + \mathfrak{B}}{t}\right)$$

arithmetic operations in the ground field where

$$\mathfrak{B} = D^s (D - 1)^{n-s} \left( 2(n - s)(D - 1) \binom{n - 1}{s - 2} + (n(D - 2) + s) \binom{n - 1}{s - 1} \right).$$

This complexity result extends the real root classification complexity of Chapter 5 to *generic determinantal systems*.

Even though our algorithm has the same complexity  $D^{O(nt)}$  as the ones based on critical point method (e.g., [9, Algo 14.6]), we make explicitly the exponent constant hidden by the above big- $O$  notation. Especially, we provide a degree bound for polynomials in the output. This degree bound  $\mathfrak{B}$  is observed to be sharp for generic inputs and, if  $s$  is fixed and  $D = 2$ ,  $\mathfrak{B}$  becomes polynomial in  $n$ .

We also emphasize that the other algorithms using critical point method are not implemented. The state-of-the-art software are based on the CAD and have a complexity  $(sD)^{2^{O(n+t)}}$ . Unlike those software which returns only explicit complicated formulas, our output formulas are encoded through the minors of parametric Hermite matrices, which are easy to be evaluated.

To support our theoretical claim, we report on the practical performance of our implementation comparing with quantifier elimination in Maple and Mathematica for both generic and

non-generic instances. Our algorithm allows us to solve quantifier elimination problems which are out of reach of these state-of-the-art software (up to 8 variables).

This is joint-work with M. Safey El Din.

## 6.1 Introduction

### 6.1.1 Problem statement

Let  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  with  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{y} = (y_1, \dots, y_t)$ . Given the quantified semi-algebraic formula  $\Psi(\mathbf{x}, \mathbf{y})$  below

$$\Psi(\mathbf{x}, \mathbf{y}) : \exists \mathbf{x} \in \mathbb{R}^n : f_1(\mathbf{x}, \mathbf{y}) = \dots = f_s(\mathbf{x}, \mathbf{y}) = 0,$$

the  $\mathbf{x}$  variables are called *quantified* variables and the  $\mathbf{y}$  variables are called *parameters*.

Solving a classical quantifier elimination problem consists in computing a *logically equivalent* quantifier-free semi-algebraic formula  $\Phi(\mathbf{y})$ , i.e.  $\Phi$  is a finite disjunction of conjunctions of polynomial constraints in  $\mathbb{Q}[\mathbf{y}]$  which is true if and only if the above quantified formula is true. Geometrically,  $\Phi$  describes the *projection* on the  $\mathbf{y}$ -space of the real algebraic set  $\mathcal{V}_{\mathbb{R}} \subset \mathbb{R}^{n+t}$  defined by the simultaneous vanishing of the  $f_i$ 's.

In this thesis, we aim at solving the following variant of quantifier elimination over the reals.

**Problem QE** (One block quantifier elimination). *Let  $\Psi(\mathbf{x}, \mathbf{y})$  be the quantified semi-algebraic formula defined as above and  $\pi : (\mathbf{x}, \mathbf{y}) \mapsto \mathbf{y}$ .*

*Design an algorithm to compute a quantifier-free semi-algebraic formula  $\Phi(\mathbf{y})$  such that  $\Phi(\mathbf{y})$  is almost equivalent with  $\Psi(\mathbf{x}, \mathbf{y})$ , i.e.,  $\Phi(\mathbf{y})$  defines a semi-algebraic subset of  $\mathbb{R}^t$  which is dense in the interior of  $\pi(\mathcal{V}_{\mathbb{R}})$ .*

**Example 6.1.1.** *We consider the toy example of the real algebraic set in  $\mathbb{R}^2$  defined by  $x^2 = y^3 - y^2$  (see Fig. 6.1). Its projection on the  $y$  coordinate is described by the quantifier-free formula*

$$y^3 - y \geq 0.$$

*For our variant quantifier elimination problem, an admissible output is*

$$y^3 - y > 0.$$

*The three endpoints  $y = 0$ ,  $y = -1$  and  $y = 1$  are dropped.*

Except for proving theorems, this is sufficient for many applications of quantifier elimination in engineering sciences (see, e.g., [144, 108, 109, 110]) where either the output formula only needs to define a sufficiently large subset of the  $\pi(\mathcal{V}_{\mathbb{R}})$  or is evaluated with parameter values which are subject to numerical noise.

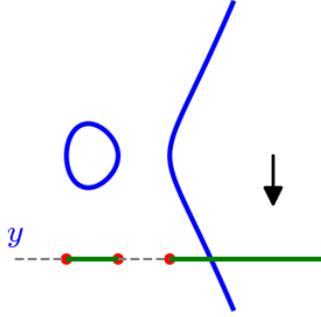


Figure 6.1: Projection of  $V(x^2 - y^3 + y)$  on  $y$ -axis.

### 6.1.2 Main results

We require that the input  $\mathbf{f} = (f_1, \dots, f_s)$  satisfies the two assumptions below.

#### Assumption 6.A.

- The ideal of  $\mathbb{Q}[\mathbf{x}, \mathbf{y}]$  generated by  $\mathbf{f}$  is radical.
- The algebraic set  $\mathcal{V} \subset \mathbb{C}^{t+n}$  of  $\mathbf{f}$  is equidimensional of dimension  $d + t$  for some  $d \in \mathbb{N}$ . Its singular locus has dimension at most  $t - 1$ .

**Assumption 6.B.** The Zariski closure  $\overline{\pi(\mathcal{V})}$  of  $\pi(\mathcal{V})$  is the whole parameter space  $\mathbb{C}^t$  and  $\pi(\mathcal{V}_{\mathbb{R}})$  is not of zero-measure in  $\mathbb{R}^t$ .

The first result of this chapter is a new probabilistic algorithm for solving the aforementioned variant of the quantifier elimination on such an input  $\mathbf{f}$ .

Our algorithm proceeds through two main steps as follows.

a) We compute a list of polynomial systems  $S_1, \dots, S_{d+1}$  in  $\mathbb{Q}[\mathbf{x}, \mathbf{y}]$  that satisfy

- Each  $S_i$  generates a zero-dimensional ideal in  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$ ;
- For almost every  $\eta \in \mathbb{R}^t$ ,

$$\bigcup_{i=1}^{d+1} (V(S_i(\eta, \cdot)) \cap \mathbb{R}^n)$$

is empty if and only if  $V(\mathbf{f}(\eta, \cdot)) \cap \mathbb{R}^n$  is empty.

This reduction is carried out by a parametric variant of the algorithm in [171] which actually computes at least one point per connected component of a regular real algebraic set. More specifically, this algorithm relies on the geometric result below, which is an extension of [171, Theorem 2].

Let  $\pi_i$  be the projection  $(x_1, \dots, x_n) \mapsto (x_1, \dots, x_i)$ . Recall that  $\text{GL}(n, t, \mathbb{C})$  denotes the change of variables that act only on  $\mathbf{x}$ . We prove that there exists a non-empty Zariski open subset  $\mathcal{A}$  of  $\text{GL}(n, t, \mathbb{C}) \times \mathbb{C}^n$  such that, for  $(A, \alpha) \in \text{GL}(n, t, \mathbb{Q}) \times \mathbb{Q}^n \cap \mathcal{A}$ , the following holds.

There exists a non-empty Zariski open subset  $\mathcal{Y}$  of  $\mathbb{C}^t$  such that, for  $\mathbf{y} \in \mathcal{Y} \cap \mathbb{R}^t$ , the union of the sets

$$\text{crit}(\pi_i, V(\mathbf{f}(\eta, \cdot)^A)) \cap \pi_{i-1}^{-1}(\alpha), \quad 1 \leq i \leq d+1,$$

contains finitely many points and meets all connected components of  $V(\mathbf{f}(\eta, \cdot)^A) \cap \mathbb{R}^n$ .

By considering  $\mathbb{Q}(\mathbf{y})$  as the ground field, the system  $S_i$  is taken as a defining system of  $\text{crit}(\pi_i, V(\mathbf{f}^A))$  using Jacobian criterion.

- b) For each system  $S_i$ , we use the real root classification algorithm described in Chapter 5 to compute a semi-algebraic formula  $\Phi_i$  whose zero set is dense in the interior of the projection of real solutions of  $S_i$ .

Finally, we return

$$\Phi = \bigvee_{i=1}^{d+1} \Phi_i$$

as the final output of the one block quantifier elimination.

A similar outline is also presented in [199, 54], in which the author computes an expensive comprehensive Gröbner systems [198] to analyze all cases before applying the real root counting algorithm of [160]. In these algorithms, the computation of comprehensive Gröbner systems and the reduction to dimension zero are known to be impractical. Whereas, our algorithm relies on the regularity assumptions of the input and the relaxation of the output to reduce to dimension zero efficiently through Safe El Din - Schost algorithm [171].

Our second goal is to analyze the complexity of this new algorithm. For generic inputs, we bound the degree of the outputs and establish a complexity result which depends on this bound. Our complexity result is then stated below. Recall that, for a fixed  $D \in \mathbb{N}$ ,  $\mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}$  denotes the set of all polynomials in  $\mathbb{C}[\mathbf{x}, \mathbf{y}]$  of total degree at most  $D$ .

**Theorem 6.1.2.** *Let*

$$\mathfrak{B} = D^s (D-1)^{n-s} \left( 2(n-s)(D-1) \binom{n-1}{s-2} + (n(D-2) + s) \binom{n-1}{s-1} \right).$$

*There exists a non-empty Zariski open subset  $\mathcal{F}$  of  $\mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^s$  such that, for every  $\mathbf{f} \in \mathcal{F}$ , our algorithm (Algorithm 6.1), in case of success, computes a semi-algebraic formula  $\Phi$  defining a dense subset of the interior of  $\pi(\mathcal{V}_{\mathbb{R}})$  within*

$$O \sim \left( 8^t \mathfrak{B}^{3t+2} \binom{t + \mathfrak{B}}{t} \right)$$

*arithmetic operations in  $\mathbb{Q}$  and  $\Phi$  involves only polynomials in  $\mathbb{Q}[\mathbf{y}]$  of degree at most  $\mathfrak{B}$ .*

Even though our algorithm has the same complexity  $D^{O(nt)}$  as the ones based on the critical point method (e.g., [9, Algo 14.6]), we make explicitly the exponent constant hidden by the big- $O$  notation. Especially, we provide also a degree bound for polynomials in the output. This degree bound  $\mathfrak{B}$  is observed to be sharp for generic inputs and, if  $s$  is fixed and  $D = 2$ ,  $\mathfrak{B}$  becomes polynomial in  $n$ .

Note that the other algorithms using critical point method are not implemented. The state-of-the-art software are based on the CAD and therefore have a complexity  $(sD)^{2O(n+t)}$ . Unlike those software which returns only explicit complicated formulas, our output formulas are encoded through the minors of parametric Hermite matrices, which are easy to be evaluated.

On the practical aspect, our implementation in MAPLE of this algorithm outperforms real quantifier elimination functions in MAPLE and MATHEMATICA. It allows us to solve examples, both generic and non-generic, that are out of reach of these software (up to 8 indeterminates). These timings are reported in Section 6.5. The degrees of polynomials involving in the output we observe from these examples agree with our theoretical bound.

**Organization of the chapter.** In Section 6.2, we recall the algorithm for real root finding of [171]. Also in the same section, we prove some auxiliary results in order to apply this algorithm parametrically. Next, we dedicate Section 6.3 for the description of our algorithm for solving the targeted problem and proving its correctness. The complexity of this algorithm is analyzed in Section 6.4. Finally, we report on some experimental results in Section 6.5.

## 6.2 Algorithm for real root finding

### 6.2.1 Safey El Din-Schost algorithm

We recall the algorithm in [171], which we refer to as the  $S^2$  algorithm, that computes at least one point per connected component of a smooth real algebraic set.

Let  $\mathbf{f} = (f_1, \dots, f_s)$  be a polynomial sequence in  $\mathbb{R}[x_1, \dots, x_n]$  that defines an algebraic set  $\mathcal{V} \subset \mathbb{C}^n$ . For  $1 \leq i \leq d$ , let  $\phi_i$  be the projection

$$\phi_i : (x_1, \dots, x_n) \mapsto (x_1, \dots, x_i).$$

We denote by  $\text{crit}(\phi_i, \mathcal{V})$  the set of critical points of the restriction of  $\phi_i$  to  $\mathcal{V}$ .

When  $\mathbf{f}$  generates a radical ideal and  $\mathcal{V}$  is a smooth equidimensional algebraic set, one can build a polynomial system using appropriate minors of  $\text{jac}(\mathbf{f})$  to define  $\text{crit}(\phi_i, \mathcal{V})$ . Note that the critical loci are nested

$$\text{crit}(\phi_1, \mathcal{V}) \subset \text{crit}(\phi_2, \mathcal{V}) \subset \dots \subset \text{crit}(\phi_d, \mathcal{V}) \subset \text{crit}(\phi_{d+1}, \mathcal{V}) = \mathcal{V}.$$

Note also that in *generic* coordinates  $\text{crit}(\phi_i, \mathcal{V})$  has expected dimension  $i - 1$  (see [171, Theorem 2]). The algorithm in [171] then exploits stronger properties of these critical loci under some genericity assumption on the coordinate system (which are satisfied through a generic linear change of coordinates).

**Proposition 6.2.1.** [171, Theorem 2] Assume that  $\mathbf{f}$  defines a smooth equidimensional algebraic set and generates a radical ideal.

Then, there exists a non-empty Zariski open set  $\mathcal{A}_{\mathbf{f}} \in \text{GL}(n, \mathbb{C})$  such that for  $A \in \mathcal{A}_{\mathbf{f}}$  the following holds:

- the restriction of  $\phi_{i-1}$  to  $\text{crit}(\phi_i, \mathcal{V}^A)$  is proper;
- the set  $\text{crit}(\phi_i, \mathcal{V}^A)$  is either empty or of dimension  $i - 1$  for  $1 \leq i \leq d + 1$ .

The first assertion in Proposition 6.2.1 implies the second one. The index in the notation  $\mathcal{A}_{\mathbf{f}}$  indicates that the non-empty Zariski open set depends on  $\mathbf{f}$ . Algorithm  $S^2$  considers fibers of the above critical loci with the convention  $\pi_0 : \mathbf{x} \rightarrow \bullet$ . Proposition 6.2.1 is the cornerstone of the  $S^2$  algorithm which can be derived from the following one.

**Proposition 6.2.2.** [171, Theorem 2] Assume that  $\mathbf{f}$  defines a smooth equidimensional algebraic set and generates a radical ideal.

Let  $\mathcal{A}_{\mathbf{f}}$  be as in Proposition 6.2.1. For  $A \in \mathcal{A}_{\mathbf{f}} \cap \text{GL}(n, \mathbb{Q})$  and  $\alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{R}^d$ , the union of the sets

$$\text{crit}(\phi_i, \mathcal{V}^A) \cap \phi_{i-1}^{-1}((\alpha_1, \dots, \alpha_{i-1})), \quad 1 \leq i \leq d + 1$$

is finite and meets all connected components of  $\mathcal{V} \cap \mathbb{R}^n$ .

**Example 6.2.3.** Let  $\mathcal{V}$  be the smooth surface defined by  $x_1^2 - x_2^2 - x_3^2 = 1$ . The Jacobian matrix  $\text{jac}(\mathbf{f})$  writes simply  $(2x_1, -2x_2, -2x_3)$ . It turns out that the identity matrix lies in the set  $\mathcal{A}$  defined in Proposition 6.2.1. Taking  $\alpha = (0, 0)$ , we obtain 3 zero-dimensional systems:

- $\text{crit}(\phi_1, \mathcal{V}) : \{-2x_2, -2x_3, x_1^2 - x_2^2 - x_3^2 - 1\}$ ,
- $\text{crit}(\phi_2, \mathcal{V}) \cap \phi_1^{-1}(\mathbf{0}) : \{-2x_3, x_1^2 - x_2^2 - x_3^2 - 1, x_1\}$ ,
- $\mathcal{V} \cap \phi_2^{-1}(\mathbf{0}) : \{x_1^2 - x_2^2 - x_3^2 - 1, x_1, x_2\}$ .

The first system admits two real solutions  $(1, 0, 0)$  and  $(-1, 0, 0)$ . The other systems do not have any real solution. The two points  $(1, 0, 0)$  and  $(-1, 0, 0)$  intersect the two connected components of  $\mathcal{V}$ .

Of course, on general examples, one would need to perform a randomly chosen linear change of variables but this example illustrates already how the algorithm works.

## 6.2.2 Parametric variant of Safey El Din - Schost algorithm

In this subsection, we describe a parametric variant of  $S^2$  algorithm. We now let  $\mathbf{f} = (f_1, \dots, f_s)$  be a polynomial sequence in  $\mathbb{Q}[\mathbf{y}][\mathbf{x}]$  where  $\mathbf{y} = (y_1, \dots, y_t)$  are considered as parameters and  $\mathbf{x} = (x_1, \dots, x_n)$  are variables. The algebraic set defined by  $\mathbf{f}$  is denoted by  $\mathcal{V} \subset \mathbb{C}^t \times \mathbb{C}^n$ . Let  $\pi$  denote the projection  $(\mathbf{x}, \mathbf{y}) \mapsto \mathbf{y}$  and  $\pi_i$  denote the projection  $(\mathbf{y}, \mathbf{x}) \mapsto (\mathbf{y}, x_1, \dots, x_i)$ .

Considering  $\mathbb{Q}(\mathbf{y})$  as the ground field, the parametric variant of  $S^2$  computes on the input  $\mathbf{f}$  a list of finite subsets of  $\mathbb{Q}[\mathbf{y}][\mathbf{x}]$ , each of which generates a zero-dimensional ideal of  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$ . These subsets are basically

$$\mathbf{f}^A \cup \Delta_i^A \cup \{x_1 - \alpha_1, \dots, x_{i-1} - \alpha_{i-1}\},$$

where  $(A, \alpha)$  is randomly chosen in  $\mathrm{GL}(n, t, \mathbb{Q}) \times \mathbb{Q}^n$  and  $\Delta_i^A$  is the set of all  $(n-d)$ -minors of the Jacobian matrix of  $\mathbf{f}^A$  with respect to  $x_i, \dots, x_n$ .

The rest of this subsection is devoted to the auxiliary results that allow us to use the  $S^2$  algorithm parametrically as above.

**Lemma 6.2.4.** *When Assumptions (6.A) and (6.B) hold, there exists a non-empty Zariski open subset  $\mathcal{B}$  of  $\mathbb{C}^t$  such that for every  $\eta \in \mathcal{B}$ , the specialization  $\mathbf{f}(\eta, \cdot)$  of  $\mathbf{f}$  at  $\eta$  generates a radical equidimensional ideal whose algebraic set is either empty or has dimension  $d$ .*

*Proof.* Under Assumption (6.B), by the fiber dimension theorem (Theorem 2.5.7), there exists a non-empty Zariski open subset  $\mathcal{B}'$  of  $\mathbb{C}^t$  such that  $\pi^{-1}(\eta) \cap \mathcal{V}$  is an algebraic set of dimension  $d$ .

Let  $\mathcal{W}$  denote the set of points of  $\mathcal{V}$  at which the Jacobian matrix  $\mathrm{jac}_{\mathbf{x}}(\mathbf{f})$  of  $\mathbf{f}$  with respect to  $\mathbf{x}$  has rank at most  $n-d-1$ . We note that  $\mathcal{W} = \mathrm{crit}(\pi, \mathcal{V}) \cup \mathrm{sing}(\mathcal{V})$ . The algebraic version of Sard's theorem [174, Proposition B2] implies that  $\pi(\mathrm{crit}(\pi, \mathcal{V}))$  is contained in a proper Zariski closed subset of  $\mathbb{C}^t$ . On the other hand, as Assumptions (6.A) hold, the dimension of  $\pi(\mathrm{sing}(\mathcal{V}))$  is less than  $t$ . Thus, it is also contained in a proper Zariski closed subset of  $\mathbb{C}^t$ .

Hence, the Zariski closure of  $\pi(\mathcal{W})$  is a proper Zariski closed subset of  $\mathbb{C}^t$ . Let  $\mathcal{B}$  be the intersection of the complement in  $\mathbb{C}^t$  of this Zariski closure with  $\mathcal{B}'$ . For  $\eta \in \mathcal{B}$ , the set

$$\{\mathbf{x} \in \mathbb{C}^n \mid \mathbf{f}(\eta, \mathbf{x}) = 0, \mathrm{rank} \mathrm{jac}_{\mathbf{x}}(\mathbf{f})(\eta) < n-d\}$$

is empty. Since the dimension of  $\pi^{-1}(\eta) \cap \mathcal{V}$  is  $d$  and the Jacobian matrix  $\mathrm{jac}_{\mathbf{x}}(\mathbf{f})(\eta, \cdot)$  of  $\mathbf{f}(\eta, \cdot)$  with respect to the variables  $\mathbf{x}$  is of rank  $n-d$  for every  $(\eta, \mathbf{x}) \in \mathcal{V} \cap \pi^{-1}(\eta)$ , the ideal  $\mathbf{f}(\eta, \cdot)$  is radical and defines a smooth and equidimensional set of dimension  $d$  by Jacobian criterion [56, Theorem 16.19].  $\square$

Lemma 6.2.4 shows that when specializing  $\mathbf{y} = (y_1, \dots, y_t)$  to a generic point  $\eta \in \mathcal{B} \cap \mathbb{R}^t$  in  $\mathbf{f}$ , one obtains  $\mathbf{f}(\eta, \cdot)$  satisfying the assumptions of Proposition 6.2.1. One could then apply Safety El Din-Schost algorithm to  $\mathbf{f}(\eta, \cdot)$  to grab sample points in the real algebraic set it defines. However, proceeding this way would lead us to use a change of variables encoded by a matrix  $A$  depending on  $\eta$ . The result below shows that choosing one generic change of variables will be valid for most of parameters' values.

**Proposition 6.2.5.** *Assume that Assumptions (6.A) and (6.B) hold. There exists a dense Zariski open subset  $\mathcal{O}$  of  $\mathrm{GL}(n, t, \mathbb{C})$  such that for every  $A \in \mathcal{O} \cap \mathrm{GL}(n, t, \mathbb{Q})$  the following holds.*

*There exists a dense Zariski open subset  $\mathcal{Y}_A$  of  $\mathbb{C}^t$  such that  $\mathcal{Y}_A$  is a subset of the Zariski open set  $\mathcal{B}$  in Lemma 6.2.4 and  $A$  lies in the Zariski open set  $\mathcal{A}_{\mathbf{f}(\eta, \cdot)}$  defined in Proposition 6.2.1 for every  $\eta \in \mathcal{Y}_A$ .*

*Proof.* Let  $\overline{\mathbb{C}(\mathbf{y})}$  denote the algebraic closure of  $\mathbb{C}(\mathbf{y})$ . We consider  $\overline{\mathbb{C}(\mathbf{y})}$  as the coefficient field. The proof of [171, Theorem 1] is purely algebraic and then is valid over the based field  $\overline{\mathbb{C}(\mathbf{y})}$ . Hence, there exists a non-empty Zariski open subset  $\tilde{\mathcal{O}}$  of  $\mathrm{GL}(n, t, \overline{\mathbb{C}(\mathbf{y})})$  such that for  $A \in \tilde{\mathcal{O}} \cap \mathrm{GL}(n, t, \mathbb{Q})$ , the variables  $(x_1, \dots, x_{i-1})$  is in Noether position with respect to the ideal in  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$  generated by  $\mathbf{f}^A + \Delta_i^A$  for  $1 \leq i \leq d+1$  where  $\Delta_i^A$  is the set of maximal minors of the truncated Jacobian matrix of  $\mathrm{jac}(\mathbf{f}^A)$  with all the partial derivatives with respect to  $\mathbf{y}$  and  $x_j$  for  $1 \leq j \leq i$  being removed (hence these minors are the ones defining  $\mathrm{crit}(\pi_i, \mathcal{V}) \cup \mathrm{sing}(\mathcal{V})$ ).

This is equivalent to the following. For  $1 \leq i \leq d+1$ ,  $i \leq j \leq n$ , there exist the polynomials  $p_{i,j} \in \mathbb{Q}(\mathbf{y})[x_1, \dots, x_{i-1}, x_j]$  such that each  $p_{i,j}$  lies in the ideal of  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$  generated by  $\mathbf{f}^A \cup \Delta_i^A$  and it is monic when considering  $x_j$  as the only variable (with the coefficients in  $\mathbb{Q}(\mathbf{y})[x_1, \dots, x_{i-1}]$ ).

The denominators of  $p_{i,j}$  are then polynomials in  $\mathbb{Q}[\mathbf{y}]$ . We choose  $\tilde{\mathcal{Y}}_A$  to be the intersection of the non-empty Zariski open set  $\mathcal{B}$  defined in Lemma 6.2.4 and the non-empty Zariski open set defined by the non-vanishing of all the denominators appeared in the  $p_{i,j}$ 's. Thus, for  $\eta \notin \tilde{\mathcal{Y}}_A$ ,  $p_{i,j}(\eta, \cdot) \in \mathbb{Q}[x_1, \dots, x_{i-1}, x_j]$  is monic in  $x_j$ . Consequently,  $(x_i, \dots, x_n)$  is in Noether position with respect to the ideal of  $\mathbb{C}[\mathbf{x}]$  generated by  $\mathbf{f}^A(\eta, \cdot) \cup \Delta_i^A(\eta, \cdot)$ . To prove the dimension of the ideal generated by  $\mathbf{f}^A \cup \Delta_i^A$ , one follows the same proof in [171] over the coefficient field  $\mathbb{Q}(\mathbf{y})$ . This leads to a non-empty Zariski set  $\mathcal{Y}_A \subset \mathbb{C}^t$  containing  $\tilde{\mathcal{Y}}_A$  such that the specializations of  $\langle \mathbf{f}^A \cup \Delta_i^A \rangle$  have expected dimension. Finally, taking  $\mathcal{O} = \tilde{\mathcal{O}} \cap \mathrm{GL}(n, t, \mathbb{C})$ , the conclusion follows.  $\square$

## 6.3 One-block quantifier elimination algorithm

### 6.3.1 Description

In this subsection, we describe our algorithm for solving our variant of the quantifier elimination problem. The input is a polynomial sequence  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  satisfying Assumptions (6.A) and (6.B). Recall that  $\pi$  denotes the projection  $(\mathbf{x}, \mathbf{y}) \mapsto \mathbf{y}$ .

Further, we denote by  $Z(\Psi)$  the set of real zeros of any quantifier-free semi-algebraic formula  $\Psi$  in the variables  $\mathbf{y}$ , i.e.,

$$Z(\Psi) = \{\mathbf{y} \in \mathbb{R}^t \mid \Psi(\mathbf{y}) \text{ is true}\}.$$

By Assumptions (6.A) and (6.B), the fiber dimension theorem [184, Theorem 1.25] implies that there exists a non-empty Zariski open subset of  $\mathbb{C}^t$  such that  $\pi^{-1}(\eta)$  has dimension  $d$ . The idea is to apply the parametric variant of Safey El Din - Schost algorithm with  $\mathbb{Q}(\mathbf{y})$  as a ground field.

More precisely, we start by picking randomly  $(A, \alpha)$  in  $\mathrm{GL}(n, t, \mathbb{Q}) \times \mathbb{Q}^n$  and apply the change of variables  $\mathbf{x} \mapsto A \cdot \mathbf{x}$  to the input  $\mathbf{f}$  to obtain a new sequence  $\mathbf{f}^A$ . As  $A$  acts only on  $\mathbf{x}$ ,  $\pi(V(\mathbf{f}^A) \cap \mathbb{R}^{n+t}) = \pi(\mathcal{V}_{\mathbb{R}})$ . Hence, a quantifier-free formula that solves our problem for  $\mathbf{f}^A$  is also a solution of the same problem for  $\mathbf{f}$ .

Let  $\mathrm{jac}_{\mathbf{x}}(\mathbf{f}^A)$  be the Jacobian matrix of  $\mathbf{f}^A$  with respect to the variables  $\mathbf{x} = (x_1, \dots, x_n)$ . The columns of  $\mathrm{jac}_{\mathbf{x}}(\mathbf{f}^A)$  is denoted by  $J_1, \dots, J_n$ . We define a subroutine  $(n-d)\mathrm{Minors}$  that

takes as input a matrix whose coefficients are in  $\mathbb{Q}[\mathbf{x}, \mathbf{y}]$  and computes all of its  $(n - d)$ -minors.

For each  $1 \leq i \leq d$ , we define the system

$$W_i^{A,\alpha} = \{\mathbf{f}^A\} \cup (n - d)\text{Minors}([J_{i+1}, \dots, J_n]) \cup \{x_1 - \alpha_1, \dots, x_{i-1} - \alpha_{i-1}\}.$$

In particular,  $W_{d+1}^{A,\alpha}$  denotes

$$\mathbf{f}^A \cup \{x_1 - \alpha_1, \dots, x_d - \alpha_d\}.$$

We prove later in Lemma 6.3.2 that, for generic  $(A, \alpha)$ , the ideals of  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$  generated by  $W_i^{A,\alpha}$  are radical and zero-dimensional.

We now solve the one block quantifier elimination problem for each of  $W_i^{A,\alpha}$ . For this step, we slightly modify Algorithm 5.4 to a subroutine called ZeroDimProjection.

This subroutine takes as input a polynomial sequence  $F \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  such that the ideal of  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$  generated by  $F$  is radical and zero-dimensional and computes a quantifier-free formula  $\Phi_F$  and a polynomial  $w_F \in \mathbb{Q}[\mathbf{y}]$  that satisfies:

- $Z(\Phi_F) \subset \pi(V(F) \cap \mathbb{R}^{n+t})$ ,
- $Z(\Phi_F) \setminus V(w_F) = \pi(V(F) \cap \mathbb{R}^{n+t}) \setminus V(w_F)$ .

Recall that Algorithm 5.4 classifies the real solutions of the system  $F$  outside a proper Zariski closed subset of the space of parameters. Its output contains a polynomial  $w_F$  defining the locus to be excluded and a list of pairs

$$\{(r_i, \Phi_i) \mid 1 \leq i \leq \ell\},$$

where  $r_i \in \mathbb{N}$  and the  $\Phi_i$ 's are quantifier-free semi-algebraic formula in  $\mathbf{y}$ .

For  $\eta \in \mathbb{R}^t$ , if  $\Phi_i(\eta)$  is true, the system  $F(\eta, \cdot)$  has exactly  $r_i$  distinct real solutions.

By simply returning the disjunction of  $\Phi_i$  corresponding to  $r_i > 0$ , we obtain the desired output for ZeroDimProjection.

Calling the subroutine ZeroDimProjection on the inputs  $W_i^{A,\alpha}$  gives us the lists of semi-algebraic formulas  $\Phi_i$ . Finally, we return

$$\Phi = \bigvee_{i=1}^{d+1} \Phi_i$$

as the output of our algorithm.

The pseudo-code in Algorithm 6.1 below summarizes our algorithm, in which the subroutine GenericDimension takes the sequence  $\mathbf{f}$  as input and computes the dimension of the ideal generated by  $\mathbf{f}$  in  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$ .

---

**Algorithm 6.1:** One-block quantifier elimination

---

**Input:** A polynomial sequence  $\mathbf{f} \in \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  satisfying Assumptions (6.A) and (6.B).

**Output:** A quantifier-free semi-algebraic formula  $\Phi$  in the variables  $\mathbf{y}$  such that  $Z(\Phi)$  is dense in the interior of  $\pi(\mathcal{V}_{\mathbb{R}})$ .

- 1 Choose randomly  $(A, \alpha) \in \text{GL}(n, \mathbb{Q}) \times \mathbb{Q}^n$
  - 2  $\mathbf{f}^A \leftarrow \mathbf{f}(A \cdot \mathbf{x})$
  - 3  $[J_1, \dots, J_n] \leftarrow \text{jac}_{\mathbf{x}}(\mathbf{f}^A)$
  - 4  $d \leftarrow \text{GenericDimension}(\mathbf{f}^A)$
  - 5 **for**  $1 \leq i \leq d+1$  **do**
  - 6      $W_i^{A, \alpha} \leftarrow \{\mathbf{f}^A\} \cup (n-d) \text{ Minors}([J_{i+1}, \dots, J_n]) \cup \{x_1 - \alpha_1, \dots, x_{i-1} - \alpha_{i-1}\}$
  - 7      $\Phi_i \leftarrow \text{ZeroDimProjection}(W_i^{A, \alpha})$
  - 8 **return**  $\Phi \leftarrow \bigvee_{i=1}^{d+1} \Phi_i$
- 

We end this subsection by an example to illustrate our algorithm.

**Example 6.3.1.** We consider the polynomial  $f = x_1^2 + y_1 x_2^2 + y_2 x_2 + y_3$  in  $\mathbb{Q}[y_1, y_2, y_3][x_1, x_2]$ . Let  $\Delta = y_2^2 - 4y_1 y_3$ . The projection of  $V(f) \cap \mathbb{R}^5$  on  $(y_1, y_2, y_3)$  is

$$(\Delta \geq 0 \wedge y_1 > 0) \vee (y_1 < 0) \vee (y_1 = 0 \wedge ((y_2 \neq 0) \vee (y_2 = 0 \wedge y_3 \leq 0))).$$

Applying the parametric variant of Safey El Din-Schost algorithm for  $A = I_3$  and  $\alpha = (0, 0)$ , we obtain 2 systems

$$W_1 = \{2y_1 x_2 + y_2, f\} \quad \text{and} \quad W_2 = \{f, x_1\}.$$

Next, we call `ZeroDimProjection` on these systems, choosing  $Q = I_2$  to simplify the calculation. We obtain then  $w_{1, \infty} = w_{2, \infty} = y_1$  and the Hermite matrices:

$$H_1 = \begin{pmatrix} 2 & 0 \\ 0 & -2y_3 + y_2^2/(2y_1) \end{pmatrix}, \quad H_2 = \begin{pmatrix} 2 & -y_2/y_1 \\ -y_2/y_1 & (-2y_1 y_3 + y_2^2)/y_1^2 \end{pmatrix}.$$

The sequences of leading principal minors are respectively  $[2, \Delta/y_1]$  and  $[2, \Delta/y_1^2]$ .

We compute then 4 points representing 4 connected components of the semi-algebraic set defined by  $y_1 \neq 0 \wedge \Delta \neq 0$ :

$$(1, 1/8, 0), (-1, 1/8, 0), (1, 1/8, 1/128), (-1, 1/8, -1/128).$$

The matrix  $H_2$  has non-zero signature over the first and second points, which both lead to the sign condition  $\Delta > 0 \wedge y_1^2 > 0$ . Thus, we have

$$\Phi_2 = (\Delta > 0 \wedge y_1^2 > 0) \wedge (y_1 \neq 0).$$

For  $H_1$ , non-zero signatures are satisfied at the first and fourth points. Evaluating the sign of  $\Delta$  and  $y_1$  at those points gives

$$\Phi_1 = ((\Delta > 0 \wedge y_1 > 0) \vee (\Delta < 0 \wedge y_1 < 0)) \wedge (y_1 \neq 0).$$

The final output is therefore  $\Phi = \Phi_1 \vee \Phi_2$ , which is equivalent to

$$\begin{aligned}\Phi &= (\Delta > 0 \wedge y_1 > 0) \vee (\Delta < 0 \wedge y_1 < 0) \vee (\Delta > 0 \wedge y_1 \neq 0) \\ &= (\Delta > 0 \wedge y_1 > 0) \vee (\Delta \neq 0 \wedge y_1 < 0).\end{aligned}$$

It is straight-forward to see that  $Z(\Phi)$  is a dense subset of  $\pi(V(f) \cap \mathbb{R}^5)$ .

### 6.3.2 Correctness

We start by proving that the polynomial sequences  $W_i^{A,\alpha}$  generate radical zero-dimensional ideals of  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$ , which is the assumption required by RealRootClassification.

**Lemma 6.3.2.** *Assume that Assumptions (6.A) and (6.B) hold. Let  $\mathcal{O}$  be the Zariski open subset of  $\mathrm{GL}(n, t, \mathbb{C})$  defined in Proposition 6.2.5 and  $A \in \mathcal{O} \cap \mathrm{GL}(n, t, \mathbb{Q})$ . There exists a non-empty Zariski open subset  $\mathcal{X}$  of  $\mathbb{C}^d$  such that for  $\alpha \in \mathcal{X} \cap \mathbb{Q}^d$ , the ideal of  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$  generated by  $W_i^{A,\alpha}$  is radical and either empty or zero-dimensional.*

*Proof.* By Proposition 6.2.5, the algebraic set defined by  $W_i^{A,\alpha}(\eta, \cdot)$  is finite when  $\eta$  varies over a non-empty Zariski open subset  $\mathcal{Y}_A$  of  $\mathbb{C}^t$ . Thus, the ideal of  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$  generated by  $W_i^{A,\alpha}$  is zero-dimensional. Now we prove that the ideal generated by  $W_i^{A,\alpha}$  is radical.

Let  $M_1^A, \dots, M_\ell^A$  be the  $(n-d)$  minors of the Jacobian matrix  $J$  associated to  $\mathbf{f}^A$  when considering only the partial derivatives with respect to  $x_{i+1}, \dots, x_n$ . Recall that  $W_i^{A,\alpha}$  is the union of  $\mathbf{f}^A$  with the  $M_1^A, \dots, M_\ell^A$  with  $x_1 - \alpha_1, \dots, x_{i-1} - \alpha_{i-1}$ . Further, we denote by  $W_i^A \subset \mathbb{Q}(\mathbf{y})[\mathbf{x}]$  the ideal generated by  $\mathbf{f}^A, M_1^A, \dots, M_\ell^A$ .

The idea is to follow [174, Definitions 3.2 and 3.3] where *charts* and *atlases* are defined for algebraic sets defined by the vanishing of  $\mathbf{f}^A$  and  $M_1^A, \dots, M_\ell^A$ .

Let  $m$  be a  $(n-d-1)$  minor of  $J$ . Without loss of generality we assume that it is the upper left such minor and let  $M_1^A, \dots, M_{d-(i-1)}^A$  be the  $(n-d)$  minors of  $J$  obtained by completing  $m$  with the  $n-d$ -th line of  $J$  and the missing column. We denote by  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]_m$  the localized ring where divisions by powers of  $m$  are allowed.

By [174, Lemma B.12] there exists a non-empty Zariski open set  $\mathcal{O}'_{m,n-d}$  such that for  $A \in \mathrm{GL}(n, t, \mathbb{C})$ , the localization of the ideal generated by  $f_1^A, \dots, f_{n-d}^A, M_1^A, \dots, M_{d-(i-1)}^A$  in the ring  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]_m$  is radical and coincides with the localization of  $W_i^A$  in  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]_m$ . By [174, Prop. 3.4], there exists a non-empty Zariski open set  $\mathcal{O}'' \subset \mathrm{GL}(n, t, \mathbb{C})$  such that for  $A \in \mathcal{O}''$ , any irreducible component of the algebraic set defined by  $W_i^A$  contains a point at which a  $(n-d-1)$  minor of  $J$  does not vanish. This implies that any primary component  $W_i^A$  whose associated algebraic set contains such a point is radical and then prime.

Now define  $\Omega$  as the intersection of  $\mathcal{O}$  (defined in Proposition 6.2.5), all non-empty Zariski open sets  $\mathcal{O}'_{m,k}$  and  $\mathcal{O}''$ . Hence, we then deduce that  $W_i^A$  generates a radical ideal. It remains to prove that there exists a non-empty Zariski open set  $\mathcal{X}_i \subset \mathbb{C}^{i-1}$  such that for  $\alpha = (\alpha_1, \dots, \alpha_{i-1}) \in \mathcal{X}_i$ ,  $\langle W_i^A \rangle + \langle x_1 - \alpha_1, \dots, x_{i-1} - \alpha_{i-1} \rangle$  is radical in  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$ . Choosing  $\alpha$

outside the set of critical values of  $\pi_i$  restricted to the algebraic set defined by  $W_i^A$  in  $\overline{\mathbb{Q}(\mathbf{y})}^n$  is enough. By Sard's theorem, this set of critical values is contained in the vanishing set of a non-zero polynomial  $\nu \in \mathbb{Q}[\mathbf{y}][\mathbf{x}]$ . Now note that it suffices to define  $\mathcal{X}_i$  as the complement of the vanishing set of the coefficients of  $\nu$  when it is seen in  $\mathbb{Q}[\mathbf{x}][\mathbf{y}]$  and  $\mathcal{X} = \bigcap_{i=1}^{d+1} \mathcal{X}_i$ .  $\square$

We prove the correctness of Algorithm 6.1 in Proposition 6.3.3 below.

**Proposition 6.3.3.** *Assume that Assumptions (6.A) and (6.B) hold. Let  $\mathcal{O} \subset \mathrm{GL}(n, t, \mathbb{C})$  and  $\mathcal{X} \subset \mathbb{C}^d$  be defined respectively in Proposition 6.2.5 and Lemma 6.3.2. Then for  $A \in \mathcal{O} \cap \mathrm{GL}(n, t, \mathbb{Q})$  and  $\alpha \in \mathcal{X} \cap \mathbb{Q}^d$ , the formula  $\Phi$  computed by Algorithm 6.1 defines a dense subset of the interior of  $\pi(\mathcal{V}_{\mathbb{R}})$ .*

*Proof.* By Lemma 6.3.2,  $W_i^{A, \alpha}$  satisfies the assumptions of RealRootClassification. Thus, the calls to RealRootClassification on  $W_i^{A, \alpha}$  are valid and return the formulas  $\Phi_i$  and the polynomials  $w_{i, \infty}$ . As  $A$  acts only on  $\mathbf{x}$ ,  $\pi(\mathcal{V}_{\mathbb{R}}^A) = \pi(\mathcal{V}_{\mathbb{R}})$ . Thus,

$$Z(\Phi_i) \subset \pi(V(W_i^{A, \alpha}) \cap \mathbb{R}^{n+t}) \subset \pi(\mathcal{V}_{\mathbb{R}}^A) = \pi(\mathcal{V}_{\mathbb{R}}).$$

Therefore,  $Z(\Phi) = \bigcup_{i=1}^{d+1} Z(\Phi_i) \subset \pi(\mathcal{V}_{\mathbb{R}})$ .

By the description of  $\Phi_i$ , for  $1 \leq i \leq d+1$ ,

$$Z(\Phi_i) \setminus V(w_{i, \infty}) = \pi(V(W_i^{A, \alpha}) \cap \mathbb{R}^{n+t}) \setminus V(w_{i, \infty}).$$

Let  $\mathcal{Y}_A$  be the non-empty Zariski open subset of  $\mathbb{C}^t$  in Proposition 6.2.5 ( $\mathcal{Y}_A$  depends on the matrix  $A$ ). We denote

$$\mathcal{W} = \bigcup_{i=1}^{d+1} V(w_{i, \infty}) \cup (\mathbb{C}^t \setminus \mathcal{Y}_A).$$

We will show that, for  $\eta \in \pi(\mathcal{V}_{\mathbb{R}}^A) \setminus \mathcal{W}$ ,  $\eta \in Z(\Phi)$ .

Since  $\eta \in \pi(\mathcal{V}_{\mathbb{R}}^A)$ ,  $V(\mathbf{f}^A(\eta, \cdot)) \cap \mathbb{R}^n$  is not empty. On the other hand, as  $\eta \in \mathcal{Y}_A$ ,  $\mathbf{f}^A(\eta, \cdot)$  generates a radical equidimensional ideal whose algebraic set is either empty or smooth of dimension  $d$ . By Proposition 6.2.2,  $V(\mathbf{f}^A(\eta, \cdot)) \cap \mathbb{R}^n$  is not empty if and only if  $\bigcup_{i=1}^{d+1} V(W_i^{A, \alpha}(\eta) \cap \mathbb{R}^n)$  is not empty either. We deduce that  $\eta \in \bigcup_{i=1}^{d+1} \pi(V(W_i^{A, \alpha}) \cap \mathbb{R}^{n+t}) \setminus \mathcal{W}$ . We have that

$$\begin{aligned} \bigcup_{i=1}^{d+1} \pi(V(W_i^{A, \alpha}) \cap \mathbb{R}^{n+t}) \setminus \mathcal{W} &= \bigcup_{i=1}^{d+1} (\pi(V(W_i^{A, \alpha}) \cap \mathbb{R}^{n+t}) \setminus \mathcal{W}) \\ &= \bigcup_{i=1}^{d+1} (Z(\Phi_i) \setminus \mathcal{W}) = (\bigcup_{i=1}^{d+1} Z(\Phi_i)) \setminus \mathcal{W}. \end{aligned}$$

Therefore,  $Z(\Phi) \setminus \mathcal{W} = \pi(\mathcal{V}_{\mathbb{R}}) \setminus \mathcal{W}$  and  $\pi(\mathcal{V}_{\mathbb{R}}) \setminus Z(\Phi)$  is of measure zero in  $\mathbb{R}^t$ . By Assumption (6.B), we conclude that  $Z(\Phi)$  is a dense subset of the interior of  $\pi(\mathcal{V}_{\mathbb{R}})$ .  $\square$

## 6.4 Complexity analysis

We now estimate the arithmetic complexity of Algorithm 6.1 once  $A \in \mathcal{O} \cap \mathrm{GL}(n, t, \mathbb{Q})$  and  $\alpha \in \mathcal{X} \cap \mathbb{Q}^d$  as in Proposition 6.2.5 are chosen randomly.

In this section, the input  $\mathbf{f}$  forms a regular sequence of  $\mathbb{Q}[\mathbf{x}, \mathbf{y}]$  (then,  $s = n - d$ ) satisfying Assumptions (6.A) and (6.B). As the calls to `RealRootClassification` on the systems  $W_i^{A,\alpha}$  are the most costly parts of our algorithm, we focus on estimating their complexities. To this end, we introduce the following assumption, which will be proved to be generic below.

**Assumption 6.C.** *Let  $F \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  and  $G$  be the reduced Gröbner basis of  $F$  with respect to the  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$  order. Then  $F$  is said to satisfy Assumption (6.C) if and only if for any  $g \in G$ , the total degree of  $g$  in both  $\mathbf{x}$  and  $\mathbf{y}$  equals the degree of  $g$  with respect to only  $\mathbf{x}$ .*

Note that, in Section 5.6, a similar assumption is introduced to establish the complexity result for solving real root classification problem on a generic input.

It is proved in Proposition 5.6.2 that, on an input  $F$  satisfying Assumption (6.C), the polynomial  $w_\infty$  in `RealRootClassification` is simply 1 and the entries of the Hermite matrix  $H_F$  are in  $\mathbb{Q}[\mathbf{y}]$ . Therefore, the `SamplePoints` subroutine is called on the sequence of leading principal minors of the parametric Hermite matrices. Then, under Assumption (6.C), a degree bound for these leading principal minors can be derived from Hilbert series of  $\mathbb{Q}(\mathbf{y})[\mathbf{x}] / \langle F \rangle$  which is explicitly known by Proposition 2.7.7 (see Lemma 5.6.3 and Lemma 5.6.4).

Following this outline, one obtains the complexity bound for `RealRootClassification` for a sequence  $F$  when Assumption 6.C holds. Finally, in Proposition 5.6.1, Assumption (6.C) is proved to hold for a generic input  $F$ .

However, the systems  $W_i^{A,\alpha}$  are not generic but equipped with a *determinantal structures* (since they are constructed using minors of some Jacobian matrices). Hence, some of the theoretical claims of Section 5.6.1 are no longer valid for these systems.

When Assumption (6.C) holds for  $W_i^{A,\alpha}$ , one can follow similar steps of Lemma 5.4.5 and Proposition 5.6.2 to bound the degree of polynomials given into `SamplePoints`. The main differences will happen in two steps:

- Can we derive explicit formulas of the degree bound from Hilbert series of  $W_i^{A,\alpha}$ ?
- For which assumptions on  $\mathbf{f}$  do the systems  $W_i^{A,\alpha}$  satisfy Assumption (6.C)?

To bypass these difficulties, we make use of nice properties of *determinantal systems*. Some notations that will be used further are introduced below.

Let  $D$  be a bound of the total degree of elements of  $\mathbf{f}$ . The zero-dimensional ideal of  $\mathbb{Q}(\mathbf{y})[\mathbf{x}]$  generated by  $W_i^{A,\alpha}$  is denoted by  $\langle W_i^{A,\alpha} \rangle$ . The quotient ring  $\mathbb{Q}(\mathbf{y})[\mathbf{x}] / \langle W_i^{A,\alpha} \rangle$  is a finite dimensional  $\mathbb{Q}(\mathbf{y})$ -vector space (Theorem 3.3.1).

Let  $G_i$  be the reduced Gröbner basis of the ideal of  $\mathbb{Q}[\mathbf{x}, \mathbf{y}]$  generated by  $W_i^{A,\alpha}$  with respect to the ordering  $\text{grevlex}(\mathbf{x}) \succ \text{grevlex}(\mathbf{y})$  and  $B_i$  be the monomial basis of  $\mathbb{Q}(\mathbf{y})[\mathbf{x}] / \langle W_i^{A,\alpha} \rangle$  constructed with  $G_i$  as explained in Section 4.4.3.

We start with the following lemma.

**Lemma 6.4.1.** *When Assumption (6.C) holds for  $W_i^{A,\alpha}$ , any leading principal minor of the matrix  $H_i$  has degree bounded by  $2 \sum_{b \in B_i} \deg(b)$ .*

*Proof.* Fixing an index  $1 \leq i \leq d + 1$ , we denote by  $\{b_{i,1}, \dots, b_{i,\delta}\}$  the elements of the basis  $B_i$  defined as above. The parametric Hermite matrix of  $S_i$  with respect to the basis  $B_i$  is

$$H_i = (h_{i,j,k})_{1 \leq j,k \leq \delta}.$$

As Assumption (6.C) holds for  $S_i$ , by Lemma 5.4.5, we deduce that the entries of  $H_i$  are elements of  $\mathbb{Q}[\mathbf{y}]$ .

Moreover, by Proposition 5.6.2, since  $S_i$  satisfies Assumption (6.C), we obtain the bound

$$\deg(h_{i,j,k}) \leq \deg(b_{i,j}) + \deg(b_{j,k}).$$

Hence, we can bound the degree of any minor of  $H_i$  by

$$2 \sum_{b \in B_i}^{\delta} \deg(b).$$

□

It remains to estimate the sum  $\sum_{b \in B_i} \deg(b)$ . A bound is obtained by simply taking the product of the highest degree appeared in  $B_i$  and its cardinality. As the Hilbert series of the quotient ring  $\mathbb{Q}(\mathbf{y})[\mathbf{x}] / \langle W_i^{A,\alpha} \rangle$  when  $\mathbf{f}$  is a generic system are known (see, e.g., [69, 187]), explicit bounds of these quantities are easily obtained.

**Lemma 6.4.2.** *Let  $B_i$  be defined as above. There exists a non-empty Zariski open subset  $\mathcal{Q}$  of  $\mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^s$  such that, for  $\mathbf{f} \in \mathcal{Q}$ , the following inequality holds for  $1 \leq i \leq d + 1$ :*

$$\sum_{b \in B_i} \deg(b) \leq D^s (D - 1)^{n-i+1-s} \left( (n - s)(D - 1) \binom{n - i}{s - 2} + \frac{1}{2}(n(D - 2) + s) \binom{n - i}{s - 1} \right).$$

*Proof.* By [61, Proposition 1], there exists a dense Zariski open subset  $\mathcal{Q}_1 \subset \mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^s$  such that for  $\mathbf{f} \in \mathcal{Q}_1$ , the Hilbert series of  $\langle W_1^{A,\alpha} \rangle$  is

$$\text{HS}_1(z) = \frac{\det(P(z^{D-1}))}{z^{(D-1)\binom{s-1}{2}}} \frac{(1 - z^D)^s (1 - z^{D-1})^{n-s}}{(1 - z)^n}$$

where  $P(z)$  is the  $(s - 1) \times (s - 1)$  matrix whose  $(i, j)$ -th entry is  $\sum_k \binom{s-i}{k} \binom{n-1-j}{k} z^k$ .

On the other hand, by [16, Corollary 14], the above Hilbert series  $\text{HS}_1(z)$  can be expressed as

$$\left( \sum_{i=0}^{s-1} \binom{n-s-1+i}{i} z^{i(D-1)} \right) (1 + z + \dots + z^{D-1})^s (1 + z + \dots + z^{D-2})^{n-s}.$$

Using similar arguments as in Lemma 5.6.4, we have that

$$\sum_{b \in B_1} \deg(b) \leq \text{HS}'(1).$$

We now aim to simplify  $\text{HS}'_1(1)$ .

Let  $A(z), B(z), C(z)$  denote respectively three factors of  $\text{HS}_1(z)$ . Then,

$$\text{HS}'_1(1) = A'(1)B(1)C(1) + A(1)B'(1)C(1) + A(1)B(1)C'(1).$$

As the  $B'(z)$  and  $C'(z)$  are simple and the formula of  $A(1)$  is well-known, we have that

$$A(1)(B'(1)C(1) + B(1)C'(1)) = \frac{1}{2} \binom{n-1}{s-1} D^s (D-1)^{n-s} (s(D-1) + (n-s)(D-2)).$$

Hence, it remains to simplify

$$A'(1) = (D-1) \sum_{i=0}^{s-1} i \binom{n-s-1-i}{i}.$$

This can be done as follows:

$$\begin{aligned} \sum_{i=0}^{s-1} (s-1-i) \binom{n-s-1-i}{i} &= \sum_{i=0}^{s-2} \sum_{j=0}^i \binom{n-s-1-j}{j} \\ &= \sum_{i=0}^{s-2} \binom{n-s-i}{i} = \binom{n-1}{s-2}. \end{aligned}$$

Therefore, we obtain

$$A'(1)B(1)C(1) = D^s (D-1)^{n-s+1} \left( (s-1) \binom{n-1}{s-1} - \binom{n-1}{s-2} \right)$$

and

$$\text{HS}'_1(1) = D^s (D-1)^{n-s} \left( (n-s)(D-1) \binom{n-1}{s-2} + \frac{1}{2} (n(D-2) + s) \binom{n-1}{s-1} \right).$$

For  $1 \leq i \leq d$ , the system  $W_i^{A,\alpha}$  can also be interpreted as the system defining the critical locus of the projection  $(x_i, \dots, x_n) \mapsto x_i$  restricted to  $V(\mathbf{f}^A(\alpha_1, \dots, \alpha_{i-1}, x_i, \dots, x_n))$ . Therefore, by replacing  $n$  by  $n-i+1$  in the above bound, we deduce that, for  $1 \leq i \leq d$ , there exists a dense Zariski open subset  $\mathcal{Q}_i \subset \mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^s$  such that

$$\sum_{b \in B_i} \deg(b) \leq D^s (D-1)^{n-i+1-s} \left( (n-s)(D-1) \binom{n-i}{s-2} + \frac{1}{2} (n(D-2) + s) \binom{n-i}{s-1} \right).$$

For  $i = d + 1$ , we can use the bound established in Lemma 5.6.4

$$\sum_{b \in B_{d+1}} \deg(b) \leq \frac{sD(D-1)^s}{2}.$$

Thus, the bound holds for  $i = d + 1$ . Taking  $\mathcal{Q} = \bigcap_{i=1}^{d+1} \mathcal{Q}_i$ , we conclude the proof.  $\square$

Further, we set

$$\mathfrak{B} = D^s(D-1)^{n-s} \left( 2(n-s)(D-1) \binom{n-1}{s-2} + (n(D-2) + s) \binom{n-1}{s-1} \right).$$

Now we show that Assumption (6.C) holds generically.

**Proposition 6.4.3.** *There exists a dense Zariski open subset  $\mathcal{P} \subset \mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^s$  such that, for every  $\mathbf{f} \in \mathcal{P}$ , there exists a dense Zariski open subset  $\mathcal{K}_{\mathbf{f}} \subset \mathrm{GL}(n, t, \mathbb{C}) \times \mathbb{C}^n$  such that for  $(A, \alpha) \in \mathcal{K}_{\mathbf{f}}$ , Assumption (6.C) holds for every system  $W_i^{A, \alpha}$ .*

*Proof.* Let  $y_{t+1}$  be a new variable and  ${}^h\mathbb{Q}[\mathbf{x}, \mathbf{y}, y_{t+1}]_D$  be the set of homogeneous polynomials in  $\mathbb{Q}[\mathbf{x}, \mathbf{y}, y_{t+1}]$  of degree  $D$ . For  $F \in \mathbb{Q}[\mathbf{x}, \mathbf{y}]$ , we denote by  ${}^hF \in \mathbb{Q}[\mathbf{x}, \mathbf{y}, y_{t+1}]$  the homogenization of  $F$  with respect to all the variables  $(\mathbf{x}, \mathbf{y})$ , that means

$${}^hF = y_{t+1}^{\deg(p)} \cdot F \left( \frac{x_1}{y_{t+1}}, \dots, \frac{x_n}{y_{t+1}}, \frac{y_1}{y_{t+1}}, \dots, \frac{y_t}{y_{t+1}} \right)$$

for each  $p \in F$ . Further,  $\langle {}^hF \rangle_h$  denotes the ideal of  $\mathbb{C}[\mathbf{x}, \mathbf{y}, y_{t+1}]$  generated by  ${}^hF$ .

We consider the following property (C1): The leading terms appearing in the reduced Gröbner basis of  $\langle {}^hF \rangle_h$  with respect to  $\mathrm{grevlex}(\mathbf{x} \succ \mathbf{y} \succ y_{t+1})$  do not involve any of the variables  $y_1, \dots, y_{t+1}$ . By the proof of Proposition 5.6.1, the property (C1) implies Assumption (6.C).

Following the proof of [7, Prop. 7], if  $y_{j+1}$  is not a zero-divisor of the quotient ring

$$\mathbb{C}[\mathbf{x}, \mathbf{y}, y_{t+1}] / \langle {}^hF, y_1, \dots, y_j \rangle_h$$

for every  $0 \leq j \leq t$ , then  $F$  satisfies the property (C1). This property means that  $(y_1, \dots, y_{t+1})$  forms a regular sequence in the quotient ring  $\mathbb{C}[\mathbf{x}, \mathbf{y}, y_{t+1}] / \langle {}^hF \rangle_h$ . We name this property as (C2).

From the proof of [187, Lemma 2.1, Lemma 2.2] and [56, Proposition 18.13], there exists a dense Zariski open subset  $\mathcal{P}_1 \subset \mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^s$  such that for  $\mathbf{f} \in \mathcal{P}_1$ , there exists a dense Zariski open subset  $\mathcal{K}_{\mathbf{f}, 1} \subset \mathrm{GL}(n, t, \mathbb{C}) \times \mathbb{C}^n$  such that for  $(A, \alpha) \in \mathcal{K}_{\mathbf{f}, 1}$ ,

- The quotient ring  $\mathbb{C}[\mathbf{x}, \mathbf{y}, y_{t+1}] / \langle {}^hW_1^{A, \alpha} \rangle_h$  is a Cohen-Macaulay ring of dimension  $t + 1$ ;
- The ideal  $\langle {}^hW_1^{A, \alpha}, y_1, \dots, y_{t+1} \rangle_h$  has dimension 0.

By the unmixedness theorem [56, Corollary 18.14],  $(y_1, \dots, y_{t+1})$  is a regular sequence over  $\mathbb{C}[\mathbf{x}, \mathbf{y}, y_{t+1}] / \langle {}^h W_1^{A,\alpha} \rangle_h$ . Thus,  $W_1^{A,\alpha}$  satisfies the property (C2) and Assumption (6.C) holds.

Similar for  $2 \leq i \leq d+1$ , we obtain dense Zariski subsets  $\mathcal{P}_i \subset \mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^s$  and  $\mathcal{K}_{f,i} \subset \text{GL}(n, t, \mathbb{C}) \times \mathbb{C}^n$  for each  $\mathbf{f} \in \mathcal{P}_i$ . Taking  $\mathcal{P} = \bigcap_{i=1}^{d+1} \mathcal{P}_i$  and  $\mathcal{K}_f = \bigcap_{i=1}^{d+1} \mathcal{K}_{f,i}$  ends the proof.  $\square$

Finally, using all the established ingredients, we finish the proof of Theorem 6.1.2 stated in the introduction.

*Proof of Theorem 6.1.2.* It is well-known that Assumptions (6.A) and (6.B) are generic. Also, the set of regular sequences is dense in  $\mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^s$ . Thus, there exists a dense Zariski open subset  $\mathcal{R} \subset \mathbb{C}[\mathbf{x}, \mathbf{y}]_{\leq D}^s$  such that for any  $\mathbf{f} \in \mathcal{R}$ ,  $\mathbf{f}$  forms a regular sequence satisfying Assumptions (6.A) and (6.B). As  $V(\mathbf{f})$  has dimension  $d+t$  and  $\mathbf{f}$  forms a regular sequence in  $\mathbb{Q}[\mathbf{x}, \mathbf{y}]$ ,  $d = n - s$ .

It remains to estimate the cost of RealRootClassification. Algorithm 6.1 consists of  $(d+1)$  calls to RealRootClassification on  $W_i^{A,\alpha}$ .

Let  $\mathcal{P}$  be the dense Zariski open set in Proposition 6.4.3 and  $\mathcal{Q} = \mathcal{P} \cap \mathcal{R}$ . Then, for  $\mathbf{f} \in \mathcal{Q}$ , SamplePoints is called on a list of polynomials in  $\mathbb{Q}[\mathbf{y}]$  of degree bounded by  $\mathfrak{B}$ . The number of principal minors is equal to the dimension of the quotient ring  $\mathbb{Q}(\mathbf{y})[\mathbf{x}] / \langle W_i^{A,\alpha} \rangle$ , which is also bounded by  $\mathfrak{B}$ . Applying Theorem 5.2.1, each call to RealRootClassification on  $W_i^{A,\alpha}$  costs at most

$$O\left(8^t \mathfrak{B}^{3t+2} \binom{t + \mathfrak{B}}{t}\right)$$

arithmetic operations in  $\mathbb{Q}$ . In total, the arithmetic complexity of Algorithm 6.1 is bounded by

$$O\left(8^t \mathfrak{B}^{3t+2} \binom{t + \mathfrak{B}}{t}\right).$$

$\square$

## 6.5 Experiments

We compare the practical behavior of Algorithm 6.1 with the commands QuantifierElimination (MAPLE's RegularChains) and Resolve (MATHEMATICA) on an Intel(R) Xeon(R) Gold 6244 3.60GHz machine of 754GB RAM. The timings are given in seconds (s.), minutes (m.) and hours (h.). The symbol  $\infty$  means that the computation is stopped after 240 hours without getting the result. We use our MAPLE implementation for Hermite matrices, in which FGB package [66] is used for Gröbner bases computation. The computation of sample points is done by RAGLIB [170] which uses MSOLVE [17] for polynomial system solving.

For RealRootClassification, we use the following notations:

- HM: timings of computing Hermite matrices and their minors.

- SP: total timings of computing the sample points.
- SIZE: the largest size of the Hermite matrices.
- DEG: the highest degree of the polynomials in the output which agrees with our theoretical bound. formulas.

We start with random dense systems. We fix the total degree  $D = 2$  and run our algorithm for various  $(t, n, s)$ . In Table 6.2, SamplePoints accounts for the major part of our timings. While our algorithm can tackle these examples, neither MAPLE nor MATHEMATICA finish within 120h. The theoretical degree bound agrees with the practical observations. This agrees with the bound given in our complexity result (Theorem 6.1.2). On smaller problems, we observe that formulas computed by MAPLE and MATHEMATICA have larger degrees than our output. Hence, these implementations, based on CAD, suffer from its doubly exponential complexity while our implementation takes advantage of the singly exponential complexity of our algorithm.

$t$	$n$	$s$	HM	SP	SIZE	DEG	MAPLE	MATHEMATICA
2	3	2	.2 s.	3 s.	8	24	$\infty$	$\infty$
2	4	2	9 s.	1 m.	12	40	$\infty$	$\infty$
2	5	2	2 m.	15 m.	16	56	$\infty$	$\infty$
2	6	2	20 m.	2.5 h.	20	72	$\infty$	$\infty$
2	7	2	1.5 h.	6 h.	24	88	$\infty$	$\infty$
3	3	2	6 s.	1 m.	8	24	$\infty$	$\infty$
3	4	2	5 m.	15 m.	12	40	$\infty$	$\infty$
3	5	2	2 h.	5 h.	16	56	$\infty$	$\infty$
3	6	2	8 h.	16 h.	20	72	$\infty$	$\infty$
4	3	2	40 s.	30 m.	8	24	$\infty$	$\infty$
4	4	2	6 h.	40 h.	12	40	$\infty$	$\infty$
5	3	2	5 m.	14 h.	8	24	$\infty$	$\infty$

Table 6.2: Generic systems with  $D = 2$

Table 6.3 shows the timings for sparse systems. Each polynomial is generated with  $D = 2$  and has  $2n$  terms. Even when Assumption (6.C) is not satisfied, our algorithm still applies. Thanks to the sparsity, the size and degree of the matrices in our algorithm are smaller than in the dense cases. Thus, our algorithm runs faster here than in Table 6.2 while these examples are out of reach of MAPLE and MATHEMATICA. We observe that the degrees of output polynomials are smaller than our theoretical bound.

Table 6.4 gives the timings for structured systems. We separate the variables  $x$  into blocks of total degree 1;  $[i, n - i]$  means that the degree in  $[x_1, \dots, x_i]$  and  $[x_{i+1}, \dots, x_n]$  are respectively

$t$	$n$	$s$	HM	SP	SIZE	DEG	MAPLE	MATHEMATICA
3	3	2	3 s.	37 s.	7	22	$\infty$	$\infty$
3	4	2	2 m.	10 m.	9	34	$\infty$	$\infty$
3	5	2	2 m.	10 m.	9	32	$\infty$	$\infty$
4	3	2	20 s.	20 m.	7	22	$\infty$	$\infty$
4	4	2	15 s.	18 m.	5	20	$\infty$	$\infty$

Table 6.3: Sparse systems with  $D = 2$

1. Here, entries of the Hermite matrices have non-trivial denominators with high degree. Computation those matrices takes the major part. However, our algorithm still outperforms the two other software. Again, the degrees of output formulas are smaller than the theoretical bound.

$t$	$n$	$s$	Block	HM	SP	SIZE	DEG	MAPLE	MATHEMATICA
3	3	2	[1, 2]	5 s.	45 s.	4	20	$\infty$	$\infty$
3	4	2	[2, 2]	4 m.	1 m.	8	32	$\infty$	$\infty$
3	5	2	[2, 3]	2 h.	9 m.	8	40	$\infty$	$\infty$
3	6	2	[3, 3]	30 h.	45 m.	14	60	$\infty$	$\infty$

Table 6.4: Structured systems

# Chapter 7

## Computing totally real hyperplane sections on algebraic curves

In this chapter, we study the following computational problem. Given an algebraic curve, embedded in projective space, we decide whether there exists a hyperplane meeting the curve in real points only. This translates into a particular type of parametrized real root counting problem that we wish to solve exactly. Combining the real root classification algorithm presented in Chapter 5 with some additional remarks, we solve a number of examples, which we can compare to the best known bounds for some classes of real algebraic curves.

Our computational method leads to the following findings:

1. There exist canonical curves  $X$  in  $\mathbb{P}^3$  with one or two ovals which do not allow any simple totally real hyperplane section (Example 7.4.1).
2. There exists a curve  $X$  in  $\mathbb{P}^3$  of genus two and degree five having one oval which does not allow any simple totally real hyperplane section (Example 7.4.2).
3. There are infinitely many plane quartics  $X$  with many ovals possessing a (complete) linear series of degree four which does not contain any totally real divisor (Example 7.5.2).

These results demonstrate the capability of our algorithm for solving applications in experimental mathematics, in particular real algebraic geometry.

This is joint-work with D. Manevich and D. Plaumann.

### 7.1 Introduction

This chapter is devoted to study the application of computing totally real hyperplane section using our real root classification described in Chapter 5.

Throughout this chapter, by a *real (algebraic) curve*  $X$ , we mean a smooth projective algebraic set of dimension 1 defined over  $\mathbb{R}$  such that the set  $X(\mathbb{R})$  of real points is non-empty (and therefore Zariski-dense in  $X$ ). Our main problem is stated as follows.

Let  $X$  be a real algebraic curve of degree  $d$  embedded into some projective space. We consider the computational problem of deciding whether there exists a real hyperplane meeting  $X$  in a prescribed number  $r$  of real points. Of particular interest is the case  $r = d$ , i.e., hyperplanes meeting  $X$  in real points only.

This problem is a special instance of the following more general problem. Given any divisor  $D$  on  $X$  defined over  $\mathbb{R}$ , and thus consisting of real points and complex-conjugate pairs, we may

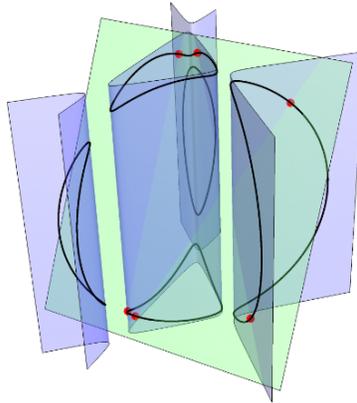


Figure 7.1: A real space curve of degree 6 with a totally real hyperplane section.

ask whether the linear series  $|D|$  (see Section 7.2 for definition) contains an effective divisor with totally real support. When  $D$  is a hyperplane section of a suitably embedded curve, it boils down to our main problem.

A number of general results have been obtained in this direction: The answer is known to be positive for any divisor of sufficiently high degree (see [123] and [177]). However, the precise degree required, relative to the genus of  $X$ , is the subject of several results and conjectures, some of which we will investigate from a computational point of view. Explicit bounds are only known if the real locus  $X(\mathbb{R})$  has many connected components (the so-called  $M$ -curves or  $(M - 1)$ -curves), by results due to Huisman [112] and Monnier [152]. On the other hand, very little is known about curves whose number of connected components is not close to maximal. Of course, the computational problem makes sense for any given curve and divisor, regardless of whether or not there is a general result covering all curves and divisors of the given kind.

A computational solution for computing a totally real hyperplane section can be achieved by classifying real roots of polynomial systems whose coefficients depend on parameters. More precisely, by considering the coefficients of the hyperplane's equation as *parameters*, one then associates a hyperplane to a point in the space of parameters. The number of real points at the intersection of the considered hyperplane with the curve may vary depending on the parameters, while the number of complex intersection points between the curve and the hyperplane is equal to the degree  $d$  for *generic* values of the parameters. (If the points are counted with intersection multiplicities and the curve is not contained in a hyperplane, this complex intersection number is equal to  $d$  for *all* values of the parameters.) Hence, from a computational point of view, we are considering a polynomial system, depending on parameters such that, when these parameters take generic values, the solution set over the complex numbers is finite.

When the input system generates a radical ideal, we use Algorithm 5.3, which is detailed in Chapter 5. Recall that this algorithm computes a finite partition of the parameter space into semi-

algebraic sets such that the number of real *simple* solutions (i.e., without multiplicities) to the input system is invariant for any value of parameters chosen in one of these sets. This allows us to derive the possible number of real roots to the input system with respect to the parameters.

From the computation, our main findings can be summarized as follows.

1. There exist canonical curves  $X$  in  $\mathbb{P}^3$  with one or two ovals which do not allow any simple totally real hyperplane section (Example 7.4.1).
2. There exists a curve  $X$  in  $\mathbb{P}^3$  of genus two and degree five having one oval which does not allow any simple totally real hyperplane section (Example 7.4.2).
3. There are infinitely many plane quartics  $X$  with many ovals possessing a (complete) linear series of degree four which does not contain any totally real divisor (Example 7.5.2).

**Organization of the chapter.** Section 7.2 is devoted to preliminaries; we recall basic definitions and properties which will be used specifically in this chapter. Section 7.3 describes how we adapt the algorithm of Chapter 5 to solve parametric polynomial systems representing the hyperplane sections. In Section 7.4, we carry out the computation on (canonical) space curves using our method. In Section 7.5, we determine the real divisor bound for certain plane quartics.

## 7.2 Preliminaries

In this section, we fix some terminology concerning real algebraic curves, divisors and linear series. As general references, we suggest [145, Chap. 7] for the theory of divisors (covering also curves defined over non-algebraically closed fields) and [100, Chap. 7] for linear series.

Recall that a *real algebraic curve*  $X$  is an integral, smooth and projective algebraic curve defined over  $\mathbb{R}$  such that the set  $X(\mathbb{R})$  of real points is non-empty. Note that a smooth curve means without any singularities, real or complex.

In particular, the set  $X(\mathbb{R})$  is an analytic manifold and decomposes into a finite number of connected components, which are called the *branches* of  $X$ . If  $X$  is embedded into the projective space  $\mathbb{P}^n$ , a branch of  $X$  is an *oval* if it meets every real hyperplane in  $\mathbb{P}^n$  in an even number of real points (counted with multiplicities), while the ones that meet hyperplanes in an odd number of points are called *pseudo-lines*. In particular, a pseudo-line has non-empty intersection with any hyperplane. By Harnack's inequality [98], we have  $s \leq g + 1$ , where  $s$  is the number of branches and  $g$  is the genus of  $X$ . The curves with  $g + 1$  (resp.  $g$ ) connected components are called  $M$ -curves (resp.  $(M - 1)$ -curves).

A totally real hyperplane section of a curve  $X$  of degree  $d$  is a hyperplane defined over  $\mathbb{R}$  that intersects  $X$  at  $d$  real points counted with multiplicities. After discussing the algorithms in Section 7.3, we will examine the following problem.

**Problem 7.2.1.** *Given a real curve  $X$  embedded in projective space, decide whether  $X$  admits a totally real hyperplane section.*

This problem arises in a more general context of real algebraic curves in which *totally real divisors* are studied.

A *divisor* on  $X$  is a formal  $\mathbb{Z}$ -linear combination of some distinct points  $P_1, \dots, P_m$  of  $X$ , i.e.,

$$D = \mu_1 P_1 + \dots + \mu_m P_m,$$

where  $\mu_i \in \mathbb{Z} \setminus \{0\}$ . The set  $\{P_1, \dots, P_m\}$  is called the *support* of  $D$ , the numbers  $\mu_1, \dots, \mu_m$  the *multiplicities* and  $\sum_{i=1}^m \mu_i$  the *degree*.

If all multiplicities in  $D$  are non-negative, the divisor  $D$  is called *effective*. If all multiplicities are equal to 1, the divisor is called *simple*.

The support of a divisor on a real curve may consist of real or complex points. However, we will only consider divisors that are *defined over*  $\mathbb{R}$  and hence conjugation-invariant, i.e., for any point in the support, its complex-conjugate appears with equal multiplicity. In particular, the non-real part of a divisor is of even degree. An effective divisor  $D$  is called *totally real* if its support consists of real points only.

For any non-zero real rational function  $f \in \mathbb{R}(X)$  on  $X$ , the divisor of zeros and poles (counted with positive or negative multiplicities, respectively) is denoted  $\text{div}(f)$ . Two divisors  $D$  and  $E$  are called *linearly equivalent* if  $E = D + \text{div}(f)$  for some  $f \in \mathbb{R}(X)^*$ . The principal divisors  $\text{div}(f)$  have degree 0, hence linear equivalence preserves the degree.

The *complete linear series associated to*  $D$ , denoted by  $|D|$ , is the set of effective divisors on  $X$  which are linearly equivalent to  $D$ . A complete linear series carries the structure of a projective space (see e.g. [100, Prop. II.7.7]). Any linear subspace for the projective space structure of a complete linear series is called a *linear series*. A linear series is called *totally real* if it contains a totally real (effective) divisor.

A base point of a given linear series is a point contained in the support of all divisors of the linear series. A linear series is called *base-point-free* if it has no base point.

For a real curve  $X$  embedded into projective space  $\mathbb{P}^n$  with degree  $d$ , any hypersurface  $Z \subset \mathbb{P}^n$  of degree  $e$  not containing  $X$  defines an effective intersection divisor  $X \cdot Z$  of degree  $de$ . The set of all intersections with hypersurfaces of a fixed degree forms a linear series on  $X$ , which may or may not be complete. Clearly, such a linear series is always base-point-free.

We are interested in determining the *real divisor bound* of a given real algebraic curve.

**Problem 7.2.2.** *Given a real curve  $X$ , determine the smallest natural number  $N(X) \in \mathbb{N}^*$  such that any divisor  $D$  of degree at least  $N(X)$  is linearly equivalent to a totally real divisor, i.e.,  $|D|$  is totally real. The number  $N(X)$  is called the real divisor bound of  $X$ .*

It was shown by Krasnov [123, Thm. 2.2] and Scheiderer [177, Cor. 2.10] that the real divisor bound is always finite. Furthermore, upper and lower bounds for  $N(X)$  were found by Huisman [112] and Monnier [152] for special classes of curves, which depend on the genus  $g$  of  $X$  only. For example, if  $X$  is an  $M$ -curve or an  $(M - 1)$ -curve, then we have

$$N(X) \leq 2g - 1.$$

However, it seems difficult to find upper bounds for curves with few branches.

An easy way to determine lower bounds for  $N(X)$  is to find a linear series with a pair of complex-conjugate base points, i.e., a non-real point that is fixed throughout the linear series. With this idea, Monnier [152, Cor. 6.2] proved the inequality

$$N(X) \geq g + 1$$

for a curve  $X$  with any number of branches.

On the other hand, when one considers only base-point-free linear series, it seems that no such lower bound is known. At the end of Section 7.5, we will construct an example of such a linear series on a plane quartic curve.

Note that according to Bertini's theorem, the generic element of a linear series on  $X$  is simple away from the base locus (see [92, Ch. 1, p. 137]). However, it may happen that a linear series contains a totally real divisor, but no simple one.

For example, the linear series of lines on the plane quartic  $X = V(x^4 + y^4 - z^4) \subset \mathbb{P}^2$  contains the totally real line section

$$X \cdot V(x - z) = 4 \cdot [1 : 0 : 1]$$

which corresponds to the intersection of  $X$  and the hyperplane  $x = z$  at  $[1 : 0 : 1]$  with multiplicity 4. On the other hand, it is easy to see that there is no simple totally real line section. This leads us to study also the *simple totally real divisors*.

**Problem 7.2.3.** *Given a real curve  $X$ , determine the smallest natural number  $N'(X) \in \mathbb{N}^*$  such that any divisor of degree at least  $N'(X)$  is linearly equivalent to a simple totally real divisor.*

We call  $N'(X)$  the *simple real divisor bound* of  $X$ . It was first introduced in [5, p. 29]. Obviously, we have  $N(X) \leq N'(X)$  and a first non-trivial result comparing  $N(X)$  and  $N'(X)$  is obtained in [5, Prop. 2.1.2], namely  $N'(X) \leq 2N(X)$ . However, it appears to be unknown if  $N(X)$  and  $N'(X)$  can ever actually be different.

One reason for the importance of the simple real divisor bound comes from the possibility of transferring results from smooth to singular curves (see [153, Thm. 4.3]). Basically, the algorithm we present in Section 7.3, which is an adapted version of the one in Chapter 5, computes simple totally real hyperplane sections. When we are mainly interested in the non-existence of totally real divisors within a linear series, i.e., in lower bounds for  $N(X)$ , this algorithm can be modified in a way explained in Section 7.3 to handle totally real hyperplane sections in general.

### 7.3 Algorithm for computing totally real hyperplane sections

Given  $(f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{x}]$  where  $\mathbf{x} = (x_1, \dots, x_n)$ , we assume that

- The sequence  $(f_1, \dots, f_s)$  generates a radical ideal in  $\mathbb{Q}[\mathbf{x}]$ .
- The affine algebraic set defined by  $f_1 = \dots = f_s = 0$  is equidimensional of dimension 1 in  $\mathbb{C}^n$  of degree  $d$ .

We consider the problem whether there exists a hyperplane with real coefficients that intersects the curve at  $d$  distinct real points; such an intersection is called a *simple* totally real hyperplane intersection.

The hyperplane is modeled by a polynomial

$$h = y_1x_1 + \cdots + y_nx_n + y_{n+1}$$

where the coefficients  $(y_1, \dots, y_{n+1})$  are considered as parameters. The problem is therefore equivalent to decide whether there exists  $(y_1, \dots, y_{n+1}) \in \mathbb{R}^{n+1}$  such that the parametric system

$$f_1 = \cdots = f_s = h = 0$$

has  $d$  distinct real solutions. Further,  $\mathbf{f}$  denotes the system  $(f_1, \dots, f_s, h)$ .

We assume that for any  $\eta \subset \mathbb{C}^{n+1}$ , the number of complex solutions to  $\mathbf{f}(\eta, \mathbf{x})$  is finite. Thus, the above problem can be solved by a real root classification of  $\mathbf{f}$ .

Since the polynomial  $h$  is homogeneous in the parameters  $\mathbf{y}$ , one can dehomogenize the system above by these successive substitutions:

$$\begin{aligned} y_1 &\rightarrow 1; \\ y_1 &\rightarrow 0, \quad y_2 \rightarrow 1; \\ &\vdots \\ y_1 &\rightarrow 0, \quad \dots, \quad y_{n-1} \rightarrow 0, \quad y_n \rightarrow 1. \end{aligned}$$

In what follows, we consider only the first substitution  $y_1 \rightarrow 1$ , so the actual parameters are  $\mathbf{y} = (y_2, \dots, y_{n+1})$ . The other computations are handled in a similar way.

In this chapter, we rely on Algorithm 5.3 to identify the possible number of real solutions of the given system. Briefly, this algorithm follows three main steps below.

- (a) We start by computing a parametric Hermite matrix  $\mathcal{H}$  associated to  $(\mathbf{f}, h) \subset \mathbb{Q}[\mathbf{y}][\mathbf{x}]$  and derive two polynomials:  $\mathbf{w}_\infty$  encoding the non-specialization locus of  $\mathcal{H}$  and  $\mathbf{w}_\mathcal{H}$  which is basically the numerator of  $\det(\mathcal{H})$ . The product  $\mathbf{w}_\infty \cdot \mathbf{w}_\mathcal{H}$  is denoted by  $\mathbf{w}$ .

The details of this step is explained in Subsection 5.4.4.

- (b) Next, we compute a set of points  $\{\mathbf{a}_1, \dots, \mathbf{a}_\ell\}$  that intersects every connected component of the semi-algebraic set of  $\mathbb{R}^n$  defined by  $\mathbf{w} \neq 0$ . This step is usually the most expensive as the polynomial  $\mathbf{w}$  may have large degree (exponential in the number of variables  $n$ ).

By Proposition 5.5.2, for any  $\eta$  varying over the connected component containing a sample point  $\mathbf{a}_i$ , the number of real solutions to  $\mathbf{f}(\eta, \mathbf{x})$  is the same as the number of real solutions to  $\mathbf{f}(\mathbf{a}_i, \mathbf{x})$ .

- (c) For  $1 \leq i \leq \ell$ , evaluate the signature of the specialized Hermite matrix  $\mathcal{H}(\mathbf{a}_i)$ , which gives the number  $r_i$  of real solutions to  $\mathbf{f}(\mathbf{a}_i, \mathbf{x})$ .

This output gives all the possible numbers of real solutions of the system  $\mathbf{f}$  over  $\mathbb{R}^n \setminus V(\mathbf{w})$ .

In most of the cases, Algorithm 5.4 is sufficient to compute a hyperplane that intersects the given curve at only real points if such a hyperplane exists. However, as the resulting classification holds only for the parameters at which  $w \neq 0$ , one still needs to investigate the vanishing locus of  $w$  to obtain a complete root classification, i.e., the number of real solutions of  $\mathbf{f}$  for every  $\eta \in \mathbb{R}^t$ .

Theoretically, this can be done using a similar routine. This consists of classifying the solutions of  $\mathbf{f}$  over the vanishing locus of  $w$ . There are several possible approaches, for instance, computing over the algebraic extension  $\mathbb{Q}[\mathbf{y}]/\langle w \rangle$  or calling the algorithm above on with  $w$  added to the input system. The first approach usually leads to high arithmetic costs while the second induces Hermite matrices of large size (depending on the degree of  $w$ ). One can also try to compute the sign conditions of the leading principal minors of  $\mathcal{H}$  while imposing a rank deficiency on the matrix. This results in deciding the emptiness of a semi-algebraic set whose defining atoms are minors of the Hermite matrix. To the best of our knowledge, these methods can be computationally difficult in practice.

Note that, for  $\eta \in V(w_{\mathcal{H}}) \setminus V(w_{\infty})$ , the system  $\mathbf{f}(\eta, \cdot)$  has less than  $d$  distinct complex solutions. Thus, if we restrict to only the simple totally real hyperplane sections, it remains only to classify the real roots for the parameters belong  $V(w_{\infty}) \cap \mathbb{R}^n$ . In the examples we consider, the polynomials  $w_{\infty}$  correspond to the hyperplanes which intersect the given curves at infinity and are factorized into polynomials of small degree (at most 3). Thus, they can be treated by calling the algorithm on the input  $\mathbf{f}$  adding each factor of  $w_{\infty}$ . Looking closer, these factors can be simplified before being sent to the above algorithm to accelerate the computation. For examples, linear factors can be handled through substitutions of variables or the quadratic factors which are sums of squares can be replaced by linear equations. Further, these processes will be explained in detail for each example.

On the contrary, handling the solutions of  $w_{\mathcal{H}}$ , where the system  $\mathbf{f}$  has multiple roots, requires expensive computations as mentioned above. Therefore, our algorithm is limited at the moment to computing simple totally real hyperplane sections.

In the particular case of one-parameter (see the examples in Section 7.5), we can obtain easily the complete root classification by evaluating the signs of leading principal minors of the matrix  $\mathcal{H}$  at real solutions of  $w$  using exact algorithms for real root isolation [200, 122].

We illustrate how to obtain a *complete* real root classification using our algorithm in the case of *one parameter* by the following example.

**Example 7.3.1.** *We consider the parametric system*

$$\mathbf{f} = \{x_1^2 + x_2^2 - y, x_1^2 + x_1x_2 - yx_2 + x_1 + y^2\},$$

where  $(x_1, x_2)$  are variables and  $y$  is the parameter. Using the ordering  $\text{grevlex}(x_1 \succ x_2) \succ y$ , we obtain the basis  $\{1, x_2, x_1, x_2^2\}$  for the quotient ring  $\mathbb{Q}[y][x_1, x_2]/\langle \mathbf{f} \rangle$  and the symmetric Hermite matrix associated to this basis

$$\mathcal{H} = \begin{bmatrix} 4 & -y - 1 & y - 1 & 2y^2 + 5y \\ * & 2y^2 + 5y & -3y^2 - y + 1 & y^3/2 - 6y^2 - 3y + 1/2 \\ * & * & -2y^2 - y & 7y^3/2 + 4y^2 - y - 1/2 \\ * & * & * & -5y^4/2 + 5y^3 + 23y^2/2 + y - 1/2 \end{bmatrix}.$$

The non-specialization polynomial  $w_\infty$  in this example is identically 1. The determinant of this Hermite matrix is

$$w = w_{\mathcal{H}} = 41y^8 + 43y^7 - 59y^6 - 204y^5 - 60y^4 + 20y^3 + 4y^2 - y.$$

This polynomial has two real solutions: 0 and  $\tilde{y} \approx 1.714$ . So, the semi-algebraic set defined by  $w \neq 0$  has three connected components and the number of distinct real solutions of  $\mathbf{f}$  is invariant over each of those connected components. More precisely,

$$\begin{aligned} y < 0 : & \quad \mathbf{f} \text{ has 0 real solution,} \\ 0 < y < \tilde{y} : & \quad \mathbf{f} \text{ has 2 real solutions,} \\ \tilde{y} < y : & \quad \mathbf{f} \text{ has 0 real solution.} \end{aligned}$$

It remains to study the roots of  $\mathbf{f}$  over two real roots of  $w$ .

For  $y = 0$ , we specialize  $y$  to 0 in the leading principal minors of the matrix  $\mathcal{H}$  and obtain the sign pattern  $(1, -1, -1, 0)$ . Thus, the system has three distinct complex solutions but only one real solution when  $y = 0$ .

It is more sophisticated for  $y = \tilde{y}$  as this solution is not in  $\mathbb{Q}$ . We evaluate the signs of the leading principal minors of  $\mathcal{H}$  at  $y = \tilde{y}$  using e.g. the command `RootFinding[Isolate]` in Maple (see [200]). We obtain the sign sequence  $(1, -1, 1, 0)$  for the leading principal minors specialized at  $y = \tilde{y}$ . Hence, the system has three distinct complex solutions but no real solution. Note that evaluating numerical approximation of  $\tilde{y}$  in  $\mathcal{H}$  could also give the same sign pattern but one needs to certify this output.

## 7.4 Real algebraic curves in $\mathbb{P}^3$

In this section, we show how our algorithm is applied to compute totally real hyperplane sections for several curves in  $\mathbb{P}^3$ . The computations are carried out by our implementation of Algorithm 5.3 on a machine of Intel(R) Xeon(R) Gold 6244 3.60GHz with 754GB RAM.

In this section,  $X$  is always assumed to be a real curve and  $g$  stands for the genus of  $X$ . If  $X$  is a real rational or real elliptic curve, by [152, Prop. 3.1],  $N(X) = 1$ . Hence, we assume  $g \geq 2$ .

We first consider canonical curves: If  $X \subset \mathbb{P}^{g-1}$  is a canonical curve having  $s \geq g - 1$  branches, then the canonical linear series, which is equal to the hyperplane linear series, is totally real. Since there are no canonical curves of genus  $g \leq 2$ , the minimal examples are plane quartic curves, i.e.,  $(d - 1)(d - 2)/2 = g = 3$ . In this case, the question of whether a plane quartic curve consisting of only one oval possesses a totally real line section is related to the undulation invariant (see [162, Thm. 4.2]). In what follows, we look at canonical curves in  $\mathbb{P}^3$  with  $g = 4$ .

**Example 7.4.1.** *In this example, we consider a finite sequence of canonical curves  $X_k$  in  $\mathbb{P}^3$ ; these curves arise as complete intersections of a cubic and a quadric. Their genus is 4 and their degree is 6. In affine coordinates, we fix the real cubic polynomial*

$$f = (x + 3)(x - y - 3)(x + y - 3) - 2.$$

1. We set

$$\begin{aligned} g_5 &= x^2 + y^2 + z^2 - 100, \\ g_4 &= (x + 3)^2 + (y + 2)^2 + z^2 - 60, \\ g_3 &= x^2 + y^2 + z^2 - 50. \end{aligned}$$

Let  $X_k$  be the algebraic curve defined by the affine ideal  $I_k = \langle f, g_k \rangle$  for  $k = 3, 4, 5$ . The curve  $X_k$  has  $k$  ovals.

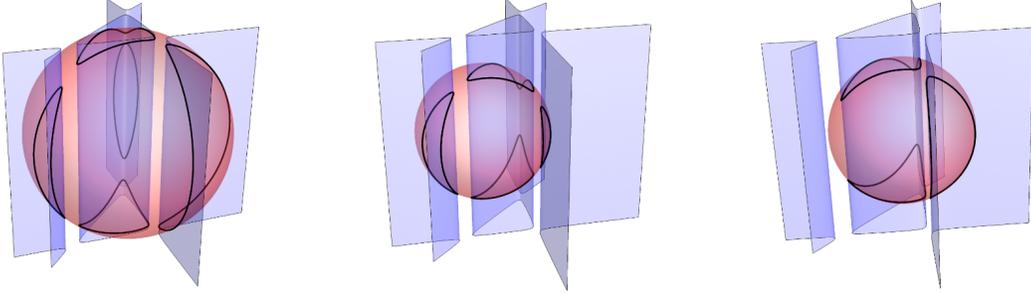


Figure 7.2: The curves  $X_5$ ,  $X_4$ , and  $X_3$ .

Computing parametric Hermite matrices on each  $I_k$  gives a boundary polynomial  $w$  of degree 18 within 5 seconds. After 10 minutes of computing of sample points for each example, we obtain affine hyperplanes which intersect the curve  $X_k$  in real points only, such as the following three hyperplanes:

$$\begin{aligned} H_5 &= x + 15307y - 8072z + 6472, \\ H_4 &= x - 14842y - 25786z - 61192, \\ H_3 &= x + 55704y - 26379z - 19751. \end{aligned}$$

Each hyperplane  $H_k$  intersects  $X_k$  in 6 (distinct) real points.

2. Setting

$$g_2 = x^2 + y^2 + z^2 - 10,$$

let  $X_2$  be the algebraic curve defined by the affine ideal  $I_2 = \langle f, g_2 \rangle$ . This curve has 2 ovals. From the theoretical point of view and in contrast to the first examples, it is a priori not clear whether this curve possesses a totally real hyperplane section. Running the algorithm for about 40 minutes on  $I_2$ , the result is that this curve does possess a totally real hyperplane section. More precisely, the hyperplane

$$H_2 = x + \frac{43}{2000}y + \frac{131}{25}z + 9,$$

intersects  $X_2$  in 6 (distinct) real points.

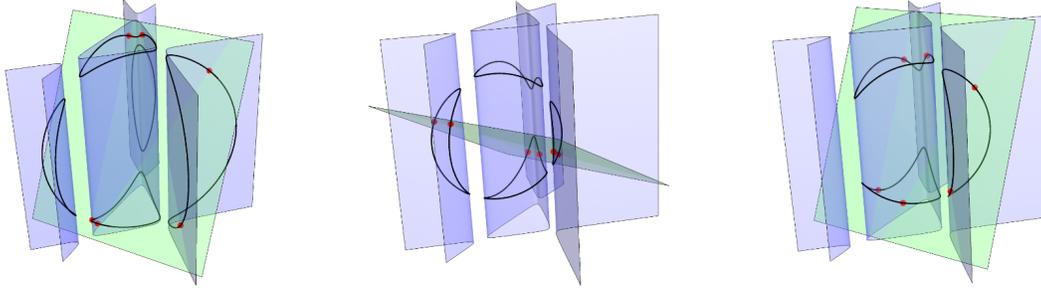


Figure 7.3: Intersecting curves and planes:  $X_i \cap H_i$  for  $i = 5, 4, 3$ .

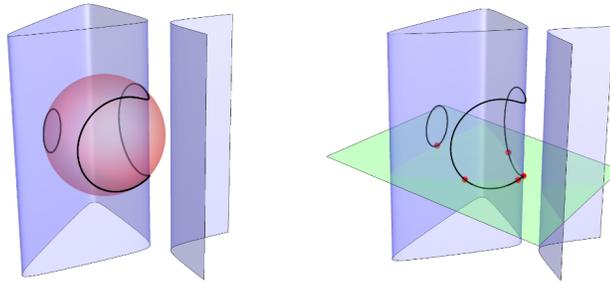


Figure 7.4: The curve  $X_2$  and its intersection with the plane  $H_2$ .

### 3. Setting

$$g'_2 = (x + 1)^2 + (y + 1)^2 + z^2 - 10,$$

let  $X'_2$  be the algebraic curve defined by the affine ideal  $I'_2 = \langle f, g'_2 \rangle$ . This curve has 2 ovals.

We compute a Hermite matrix of size  $6 \times 6$  in three parameters. From this matrix, we derive a boundary polynomial  $w$  of degree 18 with 715 monomials. These computations are done within 10 seconds.

The algorithm then computes points per connected component of the semi-algebraic set defined by  $w_\infty \cdot w_{\mathcal{H}} \neq 0$ . This computation takes almost 2 hours. In contrast to the second example, this Hermite matrix does not attain signature 6 at any of those points. Besides, the hyperplanes that correspond to the real solutions of  $w_\infty$  intersect  $X'_2$  at non-real points at infinity. Thus, these hyperplanes do not give any totally real hyperplane section.

We conclude that  $X'_2$  has no simple totally real hyperplane section. Consequently, we have  $N'(X'_2) \geq 7$ .

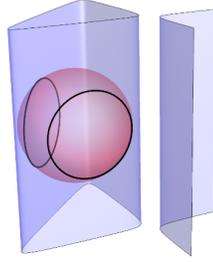


Figure 7.5: The curve  $X'_2$ .

4. For the next example, let us take the Clebsch cubic surface

$$f_0 = x^3 + y^3 + z^3 + 1 - (x + y + z + 1)^3,$$

$$g_1 = (x + 1)^2 + y^2 + z^2 - 2.$$

The algebraic curve  $X_1$  defined by the affine ideal  $I_1 = \langle f_0, g_1 \rangle$  has only 1 oval.

Our algorithm computes a  $6 \times 6$  parametric Hermite matrix in 30 seconds, from which we derive a boundary polynomial of degree 18 with 1324 monomials. Computing the sample points takes 2 hours and gives the hyperplane

$$H_1 = x - 4468y - 32932z - 10164$$

which intersects  $X_1$  in 6 (distinct) real points.

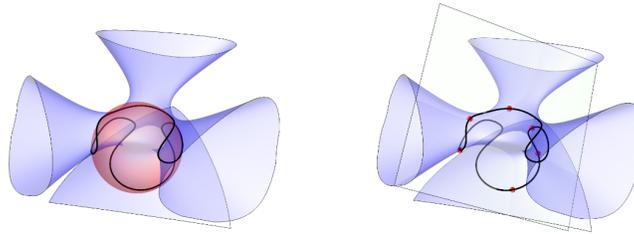


Figure 7.6: The curve  $X_1$  and its intersection with the plane  $H_1$ .

5. Finally, taking

$$g'_1 = (x + 2)^2 + y^2 + z^2 - 2,$$

let  $X_1$  be the algebraic curve defined by the affine ideal  $I'_1 = \langle f, g'_1 \rangle$ . This curve has only 1 oval, too. Again, it is a priori not clear whether this curve has a totally real hyperplane section.

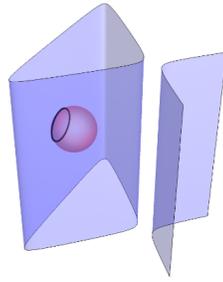


Figure 7.7: The curve  $X'_1$ .

*On this example, our algorithm behaves similarly as in the third example. We compute a  $6 \times 6$  Hermite matrix in three parameters. The boundary polynomial  $w$  has degree 18, contains 385 monomials.*

*The computation of sample points of the semi-algebraic set defined by  $w \neq 0$  takes 2 hours and none of the computed sample points gives the Hermite matrix a signature of 6. Moreover, the solutions of  $w_\infty$  here are the same as in the third example and do not correspond to a totally real hyperplane section.*

*Thus, there is no simple totally real hyperplane section in this case. Consequently, we have  $N'(X'_1) \geq 7$ .*

Of course, it takes much effort to show or disprove the existence of a canonical curve  $X$  in  $\mathbb{P}^3$  with 1 or 2 ovals and  $N(X) \leq 6$ . The existence would imply that the real divisor bound  $N(X)$  cannot depend on the main topological parameters of a real curve (the genus, the number of connected components, and whether or not the curve is of dividing type) only.

As already mentioned, it is a challenging problem to find upper bounds for  $N(X)$  in the case of curves with few branches. However, assuming a conjecture proposed by Huisman in [113, Conjecture 3.4] to be true, Monnier [152, Thm. 3.7] established new bounds for curves with  $g - 1$  connected components depending on the genus only, which is

- $N(X) \leq 3g - 1$  if  $g$  is even;
- $N(X) \leq 3g$  if  $g$  is odd.

Recently, a family of counterexamples to Huisman's conjecture has been constructed for  $n = 3$  (see [125]). These counterexamples explicitly contradict the bound found by Monnier in the case of  $g = 2$ . In the following example, we revisit two examples given in [125] through a computational approach. We will see that their construction relies on a deformation technique parameterized by a small number  $\varepsilon > 0$ . For each example, we determine different parameters  $\varepsilon$  for which there exists (and for which there does not exist) a simple totally real hyperplane section.

**Example 7.4.2.** For the first example, we consider the same polynomials as in [125, Ex. 3]. We obtain a curve of genus 4, and degree 6, which has 1 oval. In the second example, we construct a hyperelliptic curve of genus 2 and degree 5, which has 1 pseudo-line.

1. Let  $q = x_0x_3 + x_1x_2$  be the Segre quadric. We consider  $p = x_0^3 + x_1^3 + x_2^3 - x_3^3$  and

$$h = 3x_0^3 + 3x_0x_1^2 - x_0^2x_2 - 3x_0x_2^2 + x_2^3 + 4x_0^2x_3 - x_0x_1x_3 + 4x_1^2x_3 - x_2^2x_3 - 3x_0x_3^2 + x_2x_3^2 - x_3^3.$$

It is shown in [125] that the curve

$$X_\varepsilon = V(q, h + \varepsilon p)$$

does not have a totally real hyperplane section for some small parameter  $\varepsilon > 0$ . On the one hand, the algorithm shows that for  $\varepsilon = 2^{-4}$ , there is a totally real hyperplane section. For example, we can take the hyperplane  $H$  defined by

$$x_1 - \frac{323139221492926521}{1152921504606846976}x_2 + \frac{562919939027}{1099511627776}x_3 + \frac{902330031190717857}{1152921504606846976}x_0 = 0.$$

For  $\varepsilon = 2^{-5}$ , our algorithm computes a  $6 \times 6$  Hermite matrix in three parameters. The polynomial  $w_\infty$  has two factors: one is linear in the parameters and the other is a univariate polynomial of degree 3 in one parameter. The boundary polynomial  $w$  has degree 22. Computing points per connected component of the semi-algebraic set defined by  $w \neq 0$  takes about 4 hours and does not return any point that gives the Hermite matrix a signature 6.

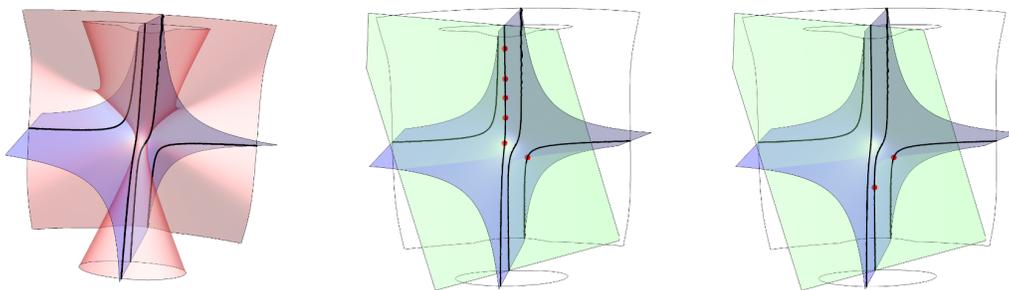


Figure 7.8: The curve  $X_{2-4}$ ; the intersections  $X_{2-4} \cap H$  and  $X_{2-5} \cap H$ .

It remains to classify the solutions when the parameters are real solutions of  $w_\infty$ . For the linear factor, we simply substitute one parameter by the others in the system to solve and use the same algorithm (with one less parameter). Finally, we call our algorithm over the algebraic extension by the univariate factor of  $w_\infty$  to classify the solutions in this case. These computations do not return any totally real hyperplane section.

So, we conclude that  $X_{2-5}$  does not have any simple totally real hyperplane section. Thus, we have  $N'(X_{2-5}) \geq 7$ .

2. In general, if  $X$  is a hyperelliptic curve, then it is known that  $N(X) \geq 2g - 1$ . If  $X$  has at least  $g$  branches, then equality holds (see [152, Cor. 6.4]).

Following [125, Cons. 1], we can construct a curve of genus 2, degree 5 with 1 pseudo-line and prescribed intersection behavior with any real hyperplane. This leads to a curve  $X_{2-8}$  below that contradicts the bound  $N(X) \leq 5$  of [152, Thm. 3.7]. Using our algorithm, we show that  $N'(X_{2-8}) \geq 6$ .

To be precise, the polynomials

$$\begin{aligned} q &= x_0x_3 - x_1x_2, \\ f &= -x_0^2x_1 - x_1^3 + 2x_0^2x_2 - x_0x_2^2 + 2x_0x_1x_3 + x_0x_2x_3 - x_0x_3^2 + x_1x_3^2, \\ g &= 2x_0x_2^2 - x_2^3 - x_0^2x_3 - x_1^2x_3 + x_2^2x_3 + 2x_0x_3^2 - x_2x_3^2 + x_3^3, \\ h_1 &= x_0^3 + x_1^3 + x_0x_2^2 - x_1x_3^2, \\ h_2 &= x_0^2x_2 + x_1^2x_3 + x_2^3 - x_3^3 \end{aligned}$$

define parametrized curves  $X_\varepsilon = V(q, f + \varepsilon h_1, g + \varepsilon h_2)$  for  $\varepsilon > 0$ . For a small parameter  $\varepsilon > 0$ , the curve  $X_\varepsilon$  does not have a totally real hyperplane section.

On the other hand, our algorithm shows that for  $\varepsilon = 2^{-4}$ , a totally real hyperplane section for  $X_{2-4}$  is given by

$$x_1 - \frac{17437072795246590045}{9223372036854775808}x_2 + \frac{8493698730591}{8796093022208}x_3 - \frac{59021162281721}{1125899906842624} = 0.$$

For  $\varepsilon = 2^{-8}$ , our algorithm computes a  $5 \times 5$  Hermite matrix in three parameters with a boundary polynomial  $w$  of degree 15. Particularly, the non-specialization polynomial  $w_\infty$  is a product of three linear polynomials of the parameters. Computing the sample points for the set defined by  $w \neq 0$  takes 3 minutes and returns no point which gives a signature 5 to the Hermite matrix.

When the parameters are real solutions of  $w_\infty$ , which has only linear factors, we substitute one parameter by the others in the parametric system. This gives us new parametric systems depending on only two parameters. Using the same algorithm, we classify the solutions of these new systems and obtain no totally real hyperplane section when  $w_\infty = 0$ . So, we conclude that there is no simple totally real hyperplane section for  $X_{2-8}$ . Thus, we have  $N'(X_{2-8}) \geq 6$ .

From the above examples, one may wonder whether it is possible to determine a largest number  $\varepsilon_0 > 0$  such that, for any  $\varepsilon \in ]0, \varepsilon_0[$ , the curve  $X_\varepsilon$  has no totally real hyperplane section. This computation can also be carried out by the algorithm we present in Section 7.3 but  $\varepsilon$  is now considered as a parameter. However, the boundary polynomial depends on 4 indeterminates and has degree up to 35. So, the computation of sample points becomes out of reach.

## 7.5 Plane quartics

Let  $X \subset \mathbb{P}^2$  be a plane quartic curve. If  $X$  has many branches, i.e., if  $s \in \{3, 4\}$ , we know that  $4 \leq N(X) \leq 5$ . We would expect  $N(X) = 5$ , so we would like to have a possibility to check if certain linear series of degree 4 do not contain a totally real divisor. The general expectation is  $N(X) = 2g - 1$  for curves of genus  $g$  having many branches (see [112, p. 92]). If  $D$  is a divisor of degree 4 on  $X$  having odd degree on at least one branch of  $X$ , then  $|D|$  can be shown to be totally real. Hence, we are interested in divisors of degree 4 having even degree on every branch. For such a divisor  $D$ , there are two possibilities. If  $D$  is special, then  $|D|$  is the canonical linear series and must be totally real. If  $D$  is non-special, then  $|D|$  defines a morphism to  $\mathbb{P}^1$  and in particular,  $D$  cannot be very ample. With the help of the algorithm, we are able to check whether each fibre of  $X \rightarrow \mathbb{P}^1$  contains a complex-conjugate pair.

If the plane quartic curve  $X$  has  $s \in \{1, 2\}$  ovals, we would like to consider very ample divisors of high degree, which give an embedding into a high-dimensional projective space. In this case, we need to check whether the hyperplane linear series of the embedded curve is totally real. For the computations, one can use the divisor package [181] in Macaulay2 [89].

**Remark 7.5.1.** *Given a plane quartic curve  $X$  with only one oval, no upper bound for  $N(X)$  is known. For two ovals, it is possible to conclude  $N(X) \leq 9$  under the assumption of an unsolved case of [113]. In particular, it is interesting to check whether every divisor of degree 10 defines a totally real linear series. If not, a new case of the conjecture is disproved. Since divisors of degree 9 on plane quartic curves are very ample, one can use the aforementioned divisor package in Macaulay2 to compute the embedding into a high-dimensional projective space. Then, one can check the existence of a totally real hyperplane section of the image curve.*

If we take a plane quartic curve  $X$  (with  $s \in \{3, 4\}$  branches) and a special divisor  $D$  of degree 4, then the linear series  $|D|$  defines a morphism  $\varphi : X \rightarrow \mathbb{P}^1$ . Using the algorithm, we can check whether there exists a real point  $[c : d] \in \mathbb{P}^1(\mathbb{R})$  which has a totally real fibre. If so, the linear series  $|D|$  is totally real. If there is no such a point, then  $|D|$  is not totally real. By perturbing the equation of the quartics (and the circles, if necessary), we get infinitely many plane quartics with many components where the real divisor bound is determined.

In what follows, we consider some examples of this type to illustrate the computation. By dehomogenizing the projective point  $[c : d]$ , we solve a polynomial system depending on one parameter. Hence, we can obtain a complete root classification of the system by the additional steps using root isolating algorithms as mentioned at the end of Section 7.3.

**Example 7.5.2.** *We continue with plane quartic curves with many branches and consider divisors of degree 4. These examples involve real root classification problems of only one parameter; our algorithm solves each of them within 2 minutes.*

1. *We can use the method described above to get a lower bound for  $N(X)$  on the curve  $X = V(x^4 + y^4 - z^4)$ . The linear series of lines is an example for a linear series which contains a*

totally real divisor, but does not contain a simple totally real one. Hence, we have  $N'(X) \geq 5$ . We consider the divisor

$$D = [1 : 0 : 1] + [0 : 1 : i] + [0 : 1 : -i] + [0 : 1 : 1]$$

which defines a morphism

$$X \rightarrow \mathbb{P}^1, \quad [x : y : z] \mapsto [xy + xz - yz - z^2 : x^2 - xz].$$

The algorithm shows that there is no totally real fibre. Even more, each fibre has of at most 2 real points. Hence, we have  $N(X) \geq 5$ .

2. In this example, we construct an explicit plane quartic curve with three ovals and a base-point-free linear series of degree four which is not totally real.

Generally, if  $X$  is a plane quartic curve and  $D$  is a special divisor of degree 4, then the morphism to  $\mathbb{P}^1$  is given by conics. Since the intersection of a quartic and a conic consists of eight points (counted with multiplicity), linear equivalence within  $|D|$  is given by a fraction of two conics having four points in common. Conversely, fixing four (real) points on  $X$ , we may consider the set of conics going through these points. The four residual points define a linear series of degree 4. Our goal is to find a linear series which is not totally real. First, we construct a plane quartic curve  $X$  with the desired topology. (There are several ways to achieve this; we use a linear determinantal representation and exploit the relation between the Cayley octad, the number of real bitangents, and the number of branches of  $X$ ; see [161]).

For example, we can take the equation of  $X$  to be

$$\begin{aligned} f = & 9x^4 - 30x^3y + 161x^2y^2 - 116xy^3 - 8y^4 + 46x^3z - 80x^2yz + 202xy^2z \\ & - 116y^3z + 59x^2z^2 - 80xyz^2 + 185y^2z^2 - 6xz^3 - 50yz^3 - 11z^4. \end{aligned}$$

Next, we take the circle  $c = x^2 + (y - \frac{z}{10})^2 - \frac{2z^2}{10}$  and fix the four real intersection points. The real vector space  $V = \text{Lin}(Q_1, Q_2)$  of conics through these points is generated by

$$\begin{aligned} Q_1 = & 0.31100521007570264x^2 - 0.4569339120067826xy \\ & + 0.7395296982938114y^2 + 0.01692042897825057xz \\ & - 0.3797243325905672yz - 0.05573253113981307z^2, \\ Q_2 = & 0.7303803360779876x^2 + 0.5870985535950933xy \\ & + 0.17978406689755905y^2 - 0.021740473005624657xz \\ & + 0.2618986086207364yz - 0.14308743118437495z^2. \end{aligned}$$

The computational problem is to check whether there is a conic in  $V$  intersecting  $X$  in only real points. As in the first example, we solve a polynomial system of one parameter using the algorithm of Section 7.3.

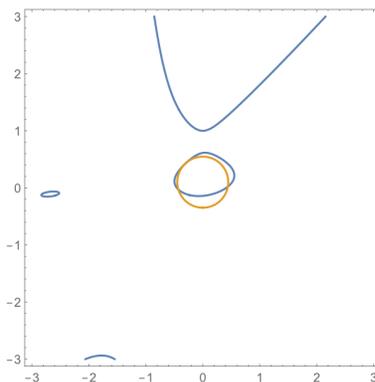


Figure 7.9: The plane quartic  $X$  and the circle  $c$ .

We start by computing a Hermite matrix of size  $8 \times 8$  and a boundary polynomial  $w$  of degree 24 ( $\deg w_\infty = 4$ ,  $\deg w_{\mathcal{H}} = 20$ ). Each fiber over the semi-algebraic set defined by  $w \neq 0$  contains 8 distinct complex points but at most 6 real points.

Next, we isolate the real solutions of  $w_{\mathcal{H}}$  and evaluate the signs of the leading principal minors of  $\mathcal{H}$  at those solutions. These sign patterns allow us to count the number of real and complex points at the real solutions of  $w_{\mathcal{H}}$ . This handles the case when the parameter takes values that satisfy  $w_{\mathcal{H}} = 0$ . For the vanishing locus of  $w_\infty$ , we call the algorithm over its associated algebraic extension. In both of these cases, we do not find any totally real fiber.

So, our algorithm shows that there is no conic in  $V$  intersecting  $X$  in real points only. Hence, taking the four residual points of any intersection  $Q \cdot X$  with  $Q \in V$  (i.e., leaving the four fixed points out), we get a divisor of degree four which does not define a totally real linear series. Furthermore, this linear series is base-point-free. The plane quartic  $X$  is an explicit example where the bound  $N(X) = 5$  is determined.

3. Analogously, we can consider the plane quartic curve  $X$  defined by

$$\begin{aligned}
 f = & (81x^4)/4 - (135x^3y)/4 + (1953x^2y^2)/16 + (297xy^3)/2 + 69y^4 \\
 & + (9x^3z)/2 + (57x^2yz)/2 + (431xy^2z)/8 - (85y^3z)/6 - (179x^2z^2)/4 \\
 & + (67xyz^2)/2 - (4685y^2z^2)/48 - (16xz^3)/3 - (1433yz^3)/36 + (917z^4)/36.
 \end{aligned}$$

The curve  $X$  consists of four ovals.

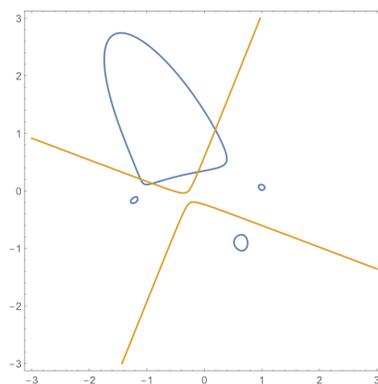


Figure 7.10: The plane quartic  $X$  and conics going through four fixed points.

*Summing up, the conics*

$$\begin{aligned}
 Q_1 &= 0.47127272928773783x^2 + 0.6598453341260914xy \\
 &\quad - 0.13447226903447518y^2 + 0.4868883263821278xz \\
 &\quad - 0.24467908024400253yz + 0.16581695886185108z^2 \\
 Q_2 &= -0.09774545786950306x^2 + 0.4442913360602867xy \\
 &\quad - 0.5056096052652832y^2 - 0.2532574091360106xz \\
 &\quad + 0.6653828276536204yz - 0.17474649814093252z^2
 \end{aligned}$$

*define the real vector space through the four fixed real points.*

*In this example, our algorithm computes a Hermite matrix  $\mathcal{H}$  of size  $8 \times 8$  and a boundary polynomial  $w$  of degree 20 ( $w_\infty = 1$ ,  $\deg w_{\mathcal{H}} = 20$ ). Again, the algorithm shows that there is no conic in this vector space intersecting  $X$  in real points only. Hence, we have  $N(X) = 5$ .*

## Chapter 8

# Computing the set of isolated points of a real algebraic set

**Abstract.** Let  $f \in \mathbb{Q}[x_1, \dots, x_n]$  be a polynomial of degree  $D$  defining an algebraic set  $\mathcal{H} \subset \mathbb{C}^n$ . We consider the problem of computing the isolated points of the real algebraic set  $\mathcal{H} \cap \mathbb{R}^n$ . This problem plays an important role for studying rigidity properties of mechanism in material designs. In this chapter, we design several algorithms for solving this computational problem.

Our algorithms share a common outline. We start with computing a finite superset of the real isolated points; the elements of this set are named the “candidates”. Such a computation is done by using the critical point method. Once the candidates are computed, it remains to identify, among them, which ones are truly real isolated points of  $\mathcal{H}$ . For this identification of isolated points, we propose two different approaches.

The first approach follows the idea of roadmaps; it constructs a real algebraic curve connecting the candidates in  $\mathcal{H} \cap \mathbb{R}^n$  whenever it is possible. The identification of isolated points boils down to answer connectivity queries in such a real algebraic curve. Using the best known bound for the complexity for computing roadmaps in [174], we obtain a probabilistic algorithm for computing the real isolated points that runs within  $(nD)^{O(n \log(n))}$  arithmetic operations in  $\mathbb{Q}$ .

The second approach decides whether a ball centered at each candidate of infinitesimal radius intersects  $\mathcal{H} \cap \mathbb{R}^n$ . However, doing this in a naive way would lead to complexity issue since the candidates are encoded by a zero-dimensional parametrization  $\mathcal{C}$  of degree  $O(D^n)$ . To bypass this difficulty, we rely on the geometric resolution algorithm over the quotient ring  $\mathbb{Q}[t]/\langle w(t) \rangle$  where  $w(t)$  is the eliminating polynomial of  $\mathcal{C}$ . This leads to an algorithm that uses  $O^\sim(64^n D^{8n})$  arithmetic operations in  $\mathbb{Q}$  and a real root isolation call on a polynomial of degree bounded by  $2^{n+2} D^{2n}$  for computing the real isolated points.

Furthermore, we also propose a variant of this algorithm to avoid partly the computations over  $\mathbb{Q}[t]/\langle w(t) \rangle$  by using “approximations” of the candidates. This variant allows us to obtain an arithmetic complexity, which lies in  $O^\sim(D^{6n+3})$  with two real root isolation calls on a polynomial of degree  $2D(D-1)^{n-1}$  and also provides a better practical performance.

Another contribution of this chapter contains several optimizations in order to achieve an efficient implementation for solving this problem. Our implementation allows us to solve instances which are out of reach of the state-of-the-art.

This chapter contains joint-works with M. Safey El Din and T. de Wolff.

## 8.1 Introduction

### 8.1.1 Problem statement

Let  $f \in \mathbb{Q}[x_1, \dots, x_n]$  and  $\mathcal{H} \subset \mathbb{C}^n$  be the hypersurface defined by  $f = 0$ . We aim at computing the *real isolated points* of  $\mathcal{H}$ , i.e. the set of points  $\mathbf{x} \in \mathcal{H} \cap \mathbb{R}^n$  such that for some positive  $r$ , the open ball centered at  $\mathbf{x}$  of radius  $r$  intersects  $\mathcal{H} \cap \mathbb{R}^n$  at only  $\mathbf{x}$ . We shall denote this set of real isolated points by  $\mathcal{I}(\mathcal{H})$ .

This problem is a particular instance of the more general one of computing the isolated points of a *semi-algebraic* set. Such problems arise naturally and frequently in the design of rigid mechanisms in material design. Those are modeled canonically with semi-algebraic constraints, and isolated points to the semi-algebraic set under consideration are related to rigidity properties of the mechanism. A particular example is the study of *auxetic* materials, i.e., materials that shrink in all directions under compression. These materials appear in nature (first discovered in [129]) e.g., in foams, bones or propylene; see e.g. [207], and have various potential applications. They are an active field of research, not only on the practical side, e.g., [95, 75], but also with respect to mathematical foundations; see e.g. [22, 23]. On the constructive side, these materials are closely related to *tensegrity frameworks*, e.g., [165, 44], which can possess various sorts of rigidity properties.

Hence, we aim to provide a practical algorithm for computing these real isolated points in the particular case of real traces of complex hypersurfaces first. This simplification allows us to significantly improve the state-of-the-art complexity for this problem and to establish a new algorithmic framework for such computations.

### 8.1.2 Main results

We provide several randomized algorithms which take as input  $f \in \mathbb{Q}[x_1, \dots, x_n]$  and compute the set of isolated points  $\mathcal{I}(\mathcal{H})$  of  $\mathcal{H} \cap \mathbb{R}^n$ . A few remarks on the data-structure are in order.

Our algorithms compute a zero-dimensional parametrization  $\mathcal{C} = (w, v_1, \dots, v_n)$  encoding a finite algebraic set

$$\mathfrak{C} = Z(\mathcal{C}) = \{(v_1(t), \dots, v_n(t)) \mid w(t) = 0\}$$

such that  $\mathfrak{C}$  contains  $\mathcal{I}(\mathcal{H})$  and a set  $\mathcal{B} = (I_1, \dots, I_s)$  of intervals in  $\mathbb{R}$  that satisfies:

- The endpoints of each interval  $I_i$  lie in  $\mathbb{Q}$ .
- Each interval  $I_i$  contains exactly one real root of  $w(t)$ , namely  $t_i$ .
- The set of isolated points of  $\mathcal{H}$  is exactly

$$\{(v_1(t_i), \dots, v_n(t_i)) \mid 1 \leq i \leq \ell\}.$$

This output represents symbolically the set of isolated points of  $\mathcal{H}$  in the sense that, for every  $\mathbf{x} \in \mathcal{I}(\mathcal{H})$ , one can derive from the pair  $(\mathcal{C}, \mathcal{B})$  a numerical representation of  $\mathbf{x}$  with any required precision.

We sketch now the geometric ingredients which allow us to compute the real isolated points of an algebraic hypersurface defined by  $f = 0$ .

Assume that  $f$  is non-negative over  $\mathbb{R}^n$  (if this is not the case, just replace it by its square) and let  $\mathbf{x} \in \mathcal{I}(\mathcal{H})$ . Since  $\mathbf{x}$  is isolated and  $f$  is non-negative over  $\mathbb{R}^n$ , the intuition is that for  $e > 0$  and small enough, the real solution set to  $f = e$  looks like a ball around  $\mathbf{x}$ , hence a bounded and closed connected component  $C_{\mathbf{x}}$ . Then  $C_{\mathbf{x}}$  contains certain critical points of the restriction of every projection on the  $x_i$ -axis to the algebraic set  $\mathcal{H}_e \subset \mathbb{C}^n$  defined by  $f = e$ . When  $e$  tends to 0, these critical points in  $C_{\mathbf{x}}$  “tend to  $\mathbf{x}$ ”. This first process allows us to compute a superset of candidate points in  $\mathcal{H} \cap \mathbb{R}^n$  containing  $\mathcal{I}(\mathcal{H})$ . Of course, one would like that this superset is finite and this will be the case up to some generic linear change of coordinates, using e.g. [169]. The elements of this finite set is called the “candidates” and are encoded by a zero-dimensional parametrization  $\mathcal{C}$ .

The next step is to identify the real isolated points of  $\mathcal{H}$  among the candidates. For this step, we follow two different strategies. The first one is based on the use of roadmaps to decide the connectivity of points over the given real algebraic set. Whereas, the second strategy aims to decide whether a ball centered of each candidate of “small” radius intersects  $\mathcal{H} \cap \mathbb{R}^n$ . In both of these strategies, we make clear the details of the algorithms.

**Constructing roadmaps.** Note that the candidates lie on “curves of critical points” which are obtained by letting  $e$  vary in the polynomial systems defining the aforementioned critical points. Assume now that  $\mathcal{H} \cap \mathbb{R}^n$  is bounded, hence contained in a ball  $B$ . Then, for  $e'$  small enough, the real algebraic set defined by  $f = 0$  is “approximated” by the union of the connected components of the real set defined by  $f = e'$  which are contained in  $B$ . Besides, these “curves of critical points”, that we just mentioned, hit these connected components when one fixes  $e'$ . We actually prove that two distinct points of our set of “candidate points” are connected through these “curves of critical points” and those connected components defined by  $f = e'$  in  $B$  if and only if they do not lie in  $\mathcal{I}(\mathcal{H})$ . Hence, we use computations of roadmaps of the real set defined by  $f = e'$  to answer those connectivity queries. Then, advanced algorithms for roadmaps and polynomial system solving allows us to achieve the announced complexity bound.

Many details are hidden in this description. In particular, we use infinitesimal deformations and techniques of semi-algebraic geometry. While infinitesimals are needed for proofs, they are difficult to use in practice. On the algorithmic side, we go further exploiting the geometry of the problem to avoid using infinitesimals.

Our complexity result for this algorithm is as follows.

**Theorem 8.1.1.** *Let  $f \in \mathbb{Q}[x_1, \dots, x_n]$  of degree  $D$  and  $\mathcal{H} \subset \mathbb{C}^n$  be the algebraic set defined by  $f = 0$ . There exists a probabilistic algorithm which, on input  $f$ , computes a zero-dimensional parametrization  $\mathcal{C}$  and isolating intervals  $\mathcal{B}$  which encode  $\mathcal{I}(\mathcal{H})$  using  $(nD)^{O(n \log(n))}$  arithmetic operations in  $\mathbb{Q}$  in case of success.*

**Solving over the quotient ring.** Recall that, at this stage, we dispose of a zero-dimensional parametrization  $\mathcal{C} = (w(t), v_1(t), \dots, v_n(t))$  that encodes candidate points. One would naturally check whether a ball of infinitesimal radius centered of each candidate intersects the algebraic set  $\mathcal{H}$ . More precisely, let  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_n) \in \mathcal{C} \cap \mathbb{R}^n$  be a candidate. We want to decide whether the system

$$f(x_1, \dots, x_n) = \sum_{i=1}^n (x_i - \eta_i)^2 - \varepsilon = 0 \quad (8.1)$$

has at least one solution in  $\mathbb{R}^n$ .

However, this direct approach faces immediately two difficulties. Writing a quantified formula to solve this decision problem raises complexity issues since those points are encoded by zero-dimensional parametrizations of degree  $D^{O(n)}$ . Besides, using an infinitesimal radius would prevent one from obtaining an efficient algorithm in practice. To bypass this difficulty, we carry out the computation over the quotient ring  $\mathbb{Q}[t]/\langle w(t) \rangle$ .

This computation relies on the results of [174, Appendix J], which provide an adaptation of the geometric resolution [87] to polynomial systems with coefficients in  $\mathbb{Q}[t]/\langle w(t) \rangle$ . Using this version of geometric resolution, the resolution of polynomial systems over  $\mathbb{Q}[t]/\langle w(t) \rangle$  induces only an additional cost of  $O^{\sim}(\deg(w))$  arithmetic operations of  $\mathbb{Q}$  comparing to the classic geometric resolution. We will see that the degree of  $w(t)$  is actually bounded by  $2D(D-1)^{n-1}$ . This allows us to obtain an algorithm that uses  $D^{O(n)}$  arithmetic operations in  $\mathbb{Q}$ .

On the algorithmic side, we go further exploiting the geometry of the problem to avoid using infinitesimals. We apply the algorithm to compute a rational number  $e_0 > 0$  that replaces the infinitesimals in the system above.

These ingredients allow us to obtain the complexity result below.

**Theorem 8.1.2.** *Let  $f \in \mathbb{Q}[x_1, \dots, x_n]$  of degree  $D$  and  $\mathcal{H} \subset \mathbb{C}^n$  be the algebraic set defined by  $f = 0$ . There exists a probabilistic algorithm which, on input  $f$ , computes a data  $(\mathcal{C}, \mathcal{B})$  encoding  $\mathcal{I}(\mathcal{H})$  in case of success using  $O^{\sim}(64^n D^{8n})$  arithmetic operations in  $\mathbb{Q}$  and one call of real root isolation on a univariate polynomial of degree bounded by  $2^{2n+2} D^{3n}$ .*

Furthermore, we propose an alternative variant that leads to a more efficient algorithm in practice. Once the rational number  $e_0$  is computed, this variant replaces the candidates by their approximations of coordinates in  $\mathbb{Q}$  and solves a similar decision problem as the one given by the system (8.1) for these approximations. Such a strategy allows us to avoid the computation over  $\mathbb{Q}[t]/\langle w(t) \rangle$ . In Subsection 8.4.5, we will define rigorously this notion of approximations.

Since this variant makes use of univariate real root isolating algorithms, a complete complexity analysis would require a bound on the bit-size of polynomials given to real root isolating algorithms. However, we observe that in practice these real root isolating steps are negligible compared to the computation over  $\mathbb{Q}[t]/\langle w(t) \rangle$ , this variant is therefore much more efficient in practice. A complexity estimate of this variant is given as below.

**Theorem 8.1.3.** *Let  $f \in \mathbb{Q}[x_1, \dots, x_n]$  of degree  $D$  and  $\mathcal{H} \subset \mathbb{C}^n$  be the algebraic set defined by  $f = 0$ . There exists a probabilistic algorithm which, on input  $f$ , computes a data  $(\mathcal{C}, \mathcal{B})$  encoding*

$\mathcal{I}(\mathcal{H})$  in case of success using  $O \sim (D^{6n+3})$  arithmetic operations in  $\mathbb{Q}$  and two real root isolating calls on a univariate polynomial of degree bounded by  $2D(D-1)^{n-1}$ .

Another contribution of this chapter consists of multiple optimization subroutines introduced for implementation. These subroutines try to avoid as much as possible the most costly computations in our algorithm. Taking into account these optimizations, we implement our algorithms in MAPLE using the libraries FGB, RAGLIB and MSOLVE. In Section 8.6, we report on practical experiments showing that they already allow us to solve non-trivial problems which are actually out of reach of [9, Alg. 12.16] which computes sample points in  $\mathcal{H} \cap \mathbb{R}^n$  only. Unfortunately, the real-life applications coming from material designs still remain intractable.

**Organization of the chapter.** In Section 8.2, we identify the set of candidates and show how to compute them. Sections 8.3 and 8.4 are respectively dedicated to present theoretical results, the descriptions and complexity analyses of our algorithms. In Section 8.5, we describe the optimizations which are used to achieve an efficient implementation. Finally, Section 8.6 reports on the practical performance of our implementation.

## 8.2 Candidates for isolated points

### 8.2.1 Identification of the candidates

This section is devoted to prove the ingredients required for computing a set of candidates. These ingredients will be used in both algorithms presented in the next sections.

As above, let  $f \in \mathbb{Q}[x_1, \dots, x_n]$  and  $\mathcal{H} \subset \mathbb{C}^n$  be the algebraic hypersurface defined by  $f = 0$ . Recall that  $\mathcal{I}(\mathcal{H})$  denotes the set of isolated points of the real algebraic set  $\mathcal{H} \cap \mathbb{R}^n$ .

**Lemma 8.2.1.** *The set  $\mathcal{I}(\mathcal{H})$  is the (finite) union of the semi-algebraically connected components of  $\mathcal{H} \cap \mathbb{R}^n$  which are a singleton.*

*Proof.* Recall that real algebraic sets have a finite number of semi-algebraically connected components [9, Theorem 5.21]. Let  $\mathcal{C}$  be a semi-algebraically connected component of  $\mathcal{H} \cap \mathbb{R}^n$ .

Assume that  $\mathcal{C}$  is not a singleton and take  $\mathbf{x}$  and  $\mathbf{y}$  in  $\mathcal{C}$  with  $\mathbf{x} \neq \mathbf{y}$ . Then, there exists a semi-algebraic continuous map  $\gamma : [0, 1] \rightarrow \mathcal{C}$  such that  $\gamma(0) = \mathbf{x}$  and  $\gamma(1) = \mathbf{y}$ . Besides, since  $\mathbf{x} \neq \mathbf{y}$ , there exist  $t \in ]0, 1[$  such that  $\gamma(t) \neq \mathbf{x}$ . By continuity of  $\gamma$  and the norm function, any ball  $B$  centered at  $\mathbf{x}$  contains  $\gamma(t) \in \mathcal{C}$  and  $\gamma(t) \neq \mathbf{x}$ . Then, any  $\mathbf{x} \in \mathcal{C}$  is not isolated in  $\mathcal{H} \cap \mathbb{R}^n$ .

Now assume that  $\mathcal{C} = \{\mathbf{x}\}$ . Observe that  $(\mathcal{H} \cap \mathbb{R}^n) \setminus \{\mathbf{x}\}$  is closed (since semi-algebraically connected components of real algebraic sets are closed). Then, the map  $\mathbf{y} \rightarrow \|\mathbf{y} - \mathbf{x}\|^2$  reaches a minimum over  $(\mathcal{H} \cap \mathbb{R}^n) \setminus \{\mathbf{x}\}$ . Let  $e$  be this minimum value. We deduce that any ball centered at  $\mathbf{x}$  of radius less than  $e$  does not meet  $(\mathcal{H} \cap \mathbb{R}^n) \setminus \{\mathbf{x}\}$ .  $\square$

To compute those connected components of  $\mathcal{H} \cap \mathbb{R}^n$  which are singletons, we use classical objects of optimization and Morse theory which are mainly *polar varieties*.

Let  $\mathbb{K}$  be an algebraically closed field, let  $\varphi \in \mathbb{K}[x_1, \dots, x_n]$  which defines the polynomial mapping

$$(x_1, \dots, x_n) \mapsto \varphi(x_1, \dots, x_n)$$

and  $V \subset \mathbb{K}^n$  be a smooth equidimensional algebraic set. We denote by  $\text{crit}(\varphi, V)$  the set of critical points of the restriction of  $\varphi$  to  $V$ . If  $c$  is the codimension of  $V$  and  $(g_1, \dots, g_s)$  generates the vanishing ideal associated to  $V$ , then  $\text{crit}(\varphi, V)$  is the subset of  $V$  at which the Jacobian matrix associated to  $(g_1, \dots, g_s, \varphi)$  has rank less than or equal to  $c$  (see Chapter 4).

In particular, the case where  $\varphi$  is replaced by the canonical projection on the  $i$ -th coordinate

$$\pi_i : (x_1, \dots, x_n) \mapsto x_i,$$

is excessively used throughout this chapter.

In our context, we do not assume that  $\mathcal{H}$  is smooth. Hence, to exploit strong topological properties of polar varieties, we retrieve a smooth situation using deformation techniques through Puiseux series (see Section 4.3).

Let  $\varepsilon$  be an infinitesimal of  $\mathbb{R}$ . By Theorem 4.3.2, the field  $\mathbb{R}\langle\varepsilon\rangle$  of Puiseux series is a real closed field and  $\mathbb{C}\langle\varepsilon\rangle = \mathbb{R}\langle\varepsilon\rangle[T]/\langle T^2 + 1 \rangle$  is its algebraic closure. We refer to Section 4.3 for the notations of initial coefficients and boundedness over  $\mathbb{R}$ .

The set of bounded elements of  $\mathbb{R}\langle\varepsilon\rangle$  is denoted by  $\mathbb{R}\langle\varepsilon\rangle_b$ . The function  $\lim_\varepsilon : \mathbb{R}\langle\varepsilon\rangle_b \rightarrow \mathbb{R}$  that maps  $\sigma$  to its *initial coefficient* is a ring homomorphism. All these definitions extend to  $\mathbb{R}\langle\varepsilon\rangle^n$  componentwise. For a semi-algebraic set  $\mathcal{S} \subset \mathbb{R}\langle\varepsilon\rangle^n$ , we naturally define the limit of  $\mathcal{S}$  as

$$\lim_\varepsilon \mathcal{S} = \left\{ \lim_\varepsilon \mathbf{x} \mid \mathbf{x} \in \mathcal{S} \text{ and } \mathbf{x} \text{ is bounded over } \mathbb{R} \right\}.$$

Given a semi-algebraic set  $\mathcal{S} \subset \mathbb{R}^n$  defined by a semi-algebraic formula  $\Phi$ ,  $\text{ext}(\mathcal{S}, \mathbb{R}\langle\varepsilon\rangle)$  denotes the (semi-algebraic) set of solutions of  $\Phi$  in  $\mathbb{R}\langle\varepsilon\rangle^n$ .

We denote by  $\mathcal{H}_\varepsilon \subset \mathbb{C}\langle\varepsilon\rangle^n$  the algebraic set defined by  $f^2 = \varepsilon^2$ . By Proposition 4.3.7,  $\mathcal{H}_\varepsilon$  is a smooth algebraic set in  $\mathbb{C}\langle\varepsilon\rangle^n$ .

Below, we give two lemmas which will be used regularly in this paper.

**Lemma 8.2.2** ([167, Lemma 3.6]). *For every  $\mathbf{x} \in \mathcal{H} \cap \mathbb{R}^n$ , there exists a point  $\mathbf{x}_\varepsilon \in \mathcal{H}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^n$  such that  $\mathbf{x}_\varepsilon$  is bounded over  $\mathbb{R}$  and  $\lim_\varepsilon \mathbf{x}_\varepsilon = \mathbf{x}$ .*

**Lemma 8.2.3** ([9, Proposition 12.51]). *Given a point  $\mathbf{x}$  lying in a bounded connected component  $\mathcal{C}$  of  $\mathcal{H} \cap \mathbb{R}^n$ , let  $\mathbf{x}_\varepsilon \in \mathcal{H}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^n$  be a point such that  $\mathbf{x}_\varepsilon$  is bounded over  $\mathbb{R}$  and  $\lim_\varepsilon \mathbf{x}_\varepsilon = \mathbf{x}$ . Let  $\mathcal{C}_\varepsilon$  be the connected component of  $\mathcal{H}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^n$  containing  $\mathbf{x}_\varepsilon$ . Then,  $\mathcal{C}_\varepsilon$  is bounded over  $\mathbb{R}$ .*

*Proof.* The proof presented in what follows are extracted from the proof of [9, Proposition 12.51]. We recall it here for the completeness of the thesis.

We prove that  $\mathcal{C}_\varepsilon$  is bounded over  $\mathbb{R}$  by contradiction. Let  $r \in \mathbb{R}$  such that  $\mathcal{C}$  is contained in the open ball  $B(\mathbf{x}, r)$ . Let  $\mathcal{C}'_\varepsilon$  be the semi-algebraically connected component of  $\mathcal{C}_\varepsilon \cap \text{ext}(B(\mathbf{x}, r), \mathbb{R}\langle\varepsilon\rangle)$ ; note that  $\mathcal{C}'_\varepsilon$  contains  $\mathbf{x}_\varepsilon$ .

We assume that  $\mathcal{C}_\varepsilon \setminus \mathcal{C}'_\varepsilon \neq \emptyset$ . Let  $z \in \mathcal{C}_\varepsilon \setminus \mathcal{C}'_\varepsilon$ , we take a semi-algebraically connected path  $\gamma : \text{ext}([0, 1], \mathbb{R}\langle\varepsilon\rangle) \rightarrow \mathcal{C}_\varepsilon$  with  $\gamma(0) = \mathbf{x}_\varepsilon$  and  $\gamma(1) = z$ . Since  $z \notin \mathcal{C}'_\varepsilon$ , the image of  $\gamma$  is not contained in  $\text{ext}(B(\mathbf{x}, r), \mathbb{R}\langle\varepsilon\rangle)$ .

Let  $t_0$  be the smallest  $t \in \text{ext}([0, 1], \mathbb{R}\langle\varepsilon\rangle)$  such that  $\gamma(t) \notin \text{ext}(B(\mathbf{x}, r), \mathbb{R}\langle\varepsilon\rangle)$  and notice that  $\mathbf{y} = \gamma(t_0) \in \text{ext}(S(\mathbf{x}, r), \mathbb{R}\langle\varepsilon\rangle)$  where  $S(\mathbf{x}, r)$  is the boundary sphere of  $B(\mathbf{x}, r)$ . Then  $\lim_\varepsilon \gamma([0, t_0])$  is connected and contained in  $\mathcal{C}$ . Hence,  $\lim_\varepsilon \mathbf{y}$  lies in  $S(\mathbf{x}, r)$  and  $\mathcal{C}$  at the same time. This contradicts  $\mathcal{C} \subset B(\mathbf{x}, r)$  and ends the proof.  $\square$

The following proposition allows us to obtain a subset of  $\mathcal{H}$  that contains  $\mathcal{S}(\mathcal{H})$ . Further, we denote this subset by  $\mathfrak{C}$  and call its elements the candidates.

**Proposition 8.2.4.** *Assume that  $\mathcal{S}(\mathcal{H})$  is not empty and let  $\mathbf{x} \in \mathcal{S}(\mathcal{H})$ . There exists a semi-algebraically connected component  $\mathcal{C}_\varepsilon$  that is bounded over  $\mathbb{R}$  of  $\mathcal{H}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^n$  such that  $\lim_\varepsilon \mathcal{C}_\varepsilon = \{\mathbf{x}\}$ .*

*Consequently, for  $1 \leq i \leq n$ , there exists an  $\mathbf{x}_\varepsilon \in \text{crit}(\pi_i, \mathcal{H}_\varepsilon) \cap \mathcal{C}_\varepsilon$  such that  $\lim_\varepsilon \mathbf{x}_\varepsilon = \mathbf{x}$ . Let  $\mathfrak{C} := \bigcap_{i=1}^n \lim_\varepsilon \text{crit}(\pi_i, \mathcal{H}_\varepsilon)$ . Then, we have*

$$\mathcal{S}(\mathcal{H}) \subset \mathfrak{C} \cap \mathbb{R}^n.$$

*Proof.* By Lemma 8.2.2, there exists  $\mathbf{x}_\varepsilon \in \mathcal{H}_\varepsilon$  such that  $\lim_\varepsilon \mathbf{x}_\varepsilon = \mathbf{x}$ . Assume that  $\mathbf{x}_\varepsilon \in \mathcal{H}_\varepsilon$  and let  $\mathcal{C}_\varepsilon$  be the connected component of  $\mathcal{H}_\varepsilon$  containing  $\mathbf{x}_\varepsilon$ . Again, by Lemma 8.2.3,  $\mathcal{C}_\varepsilon$  is bounded over  $\mathbb{R}$ . We prove that  $\lim_\varepsilon \mathcal{C}_\varepsilon = \{\mathbf{x}\}$  by contradiction.

Assume that there exists a point  $\mathbf{y}_\varepsilon \in \mathcal{C}_\varepsilon$  such that  $\lim_\varepsilon \mathbf{y}_\varepsilon = \mathbf{y}$  and  $\mathbf{y} \neq \mathbf{x}$ . Since  $\mathcal{C}_\varepsilon$  is semi-algebraically connected, there exists a semi-algebraically continuous function

$$\gamma : \text{ext}([0, 1], \mathbb{R}\langle\varepsilon\rangle) \rightarrow \mathcal{C}_\varepsilon$$

such that  $\gamma(0) = \mathbf{x}_\varepsilon$  and  $\gamma(1) = \mathbf{y}_\varepsilon$ . By [9, Proposition 12.49],  $\lim_\varepsilon \text{Im}(\gamma)$  is connected and contains  $\mathbf{x}$  and  $\mathbf{y}$ . As  $\lim_\varepsilon$  is a ring homomorphism,  $f(\lim_\varepsilon \gamma(t)) = \lim_\varepsilon f(\gamma(t)) = 0$ , so  $\lim_\varepsilon \text{Im}(\gamma)$  is contained in  $\mathcal{H} \cap \mathbb{R}^n$ . This contradicts the isolatedness of  $\mathbf{x}$ , then we conclude that  $\lim_\varepsilon \mathcal{C}_\varepsilon = \{\mathbf{x}\}$ .

Since  $\mathcal{C}_\varepsilon$  is a semi-algebraically connected component of the real algebraic set  $\mathcal{H}_\varepsilon$ , it is closed. Also,  $\mathcal{C}_\varepsilon$  is bounded over  $\mathbb{R}$ . Hence, for any  $1 \leq i \leq n$ , the projection  $\pi_i$  reaches its extrema over  $\mathcal{C}_\varepsilon$  [9, Proposition 7.6], which implies that  $\mathcal{C}_\varepsilon \cap \text{crit}(\pi_i, \mathcal{V}_\varepsilon)$  is non-empty. Take  $\mathbf{x}_\varepsilon \in \text{crit}(\pi_i, \mathcal{H}_\varepsilon) \cap \mathcal{C}_\varepsilon$ , then  $\mathbf{x}_\varepsilon$  is bounded over  $\mathbb{R}$  and its limit is  $\mathbf{x}$ . Thus,  $\mathcal{S}(\mathcal{H}) \subset \lim_\varepsilon \text{crit}(\pi_i, \mathcal{V}_\varepsilon) \cap \mathbb{R}\langle\varepsilon\rangle_b^n$  for any  $1 \leq i \leq n$ , which implies  $\mathcal{S}(\mathcal{H}) \subset \bigcap_{i=1}^n \lim_\varepsilon \text{crit}(\pi_i, \mathcal{V}_\varepsilon) \cap \mathbb{R}\langle\varepsilon\rangle_b^n$ .  $\square$

## 8.2.2 Computation of candidates

This subsection is devoted to describe a subroutine Candidates to compute a zero-dimensional parametrization  $\mathcal{C}$  encoding the candidates defined in Proposition 8.2.4.

By Proposition 8.2.4, the set

$$\mathfrak{C} = \bigcap_{i=1}^n \lim_\varepsilon \text{crit}(\pi_i, \mathcal{H}_\varepsilon)$$

contains the real isolated points of  $\mathcal{H}$ . To ensure that this set is finite, we use a *generically chosen* linear change of coordinates.

Given a matrix  $A \in \text{GL}(n, \mathbb{Q})$ , let  $f^A = f(A \cdot \mathbf{x})$  and  $\mathcal{H}^A = V(f^A) \subset \mathbb{C}^n$ . The algebraic subset of  $\mathbb{C}\langle \varepsilon \rangle^n$  defined by  $(f^A - \varepsilon) \cdot (f^A + \varepsilon) = 0$  is denoted by  $\mathcal{H}_\varepsilon^A$ .

Let  $y$  be a new variable. For  $1 \leq i \leq n$ ,  $I_i$  denotes the ideal of  $\mathbb{Q}[y, x_1, \dots, x_n]$  generated by the set of polynomials

$$\left\{ y \cdot \frac{\partial f^A}{\partial x_i} - 1, \frac{\partial f^A}{\partial x_j} \quad \text{for all } j \neq i \right\}.$$

Our subroutine *Candidates* relies on the geometric results presented in [171] and [169]. In [171], it is proved that every  $\text{crit}(\pi_i, \mathcal{H}_\varepsilon^A)$  is a finite set when  $A$  is taken outside a prescribed proper Zariski closed subset of  $\text{GL}(n, \mathbb{C})$ .

Moreover, as  $A$  is assumed to be generically chosen, [169, Theorem 1] shows that the algebraic set associated to the ideal

$$\langle f^A \rangle + (I_i \cap \mathbb{Q}[x_1, \dots, x_n])$$

is finite and contains  $\lim_\varepsilon \text{crit}(\pi_i, \mathcal{H}_\varepsilon^A)$ .

Note that, for any matrix  $A$ , the real isolated points of  $\mathcal{H}^A$  is the image of  $\mathcal{S}(\mathcal{H})$  by the linear mapping associated to  $A^{-1}$ . Thus, in practice, we will choose randomly a matrix  $A \in \text{GL}(n, \mathbb{Q})$ , compute the real isolated points of  $\mathcal{H}^A$ , and then go back to  $\mathcal{S}(\mathcal{H})$  by applying the change of coordinates induced by  $A^{-1}$ . This random choice of  $A$  makes the subroutine *Candidates* probabilistic.

In our problem, the intersection of  $\lim_\varepsilon \text{crit}(\pi_i, \mathcal{H}_\varepsilon^A)$  is needed rather than each limit itself. Hence, we use the inclusion

$$\mathcal{S}(\mathcal{H}^A) \subset \bigcap_{i=1}^n \lim_\varepsilon \text{crit}(\pi_i, \mathcal{H}_\varepsilon^A) \subset V \left( \langle f^A \rangle + \sum_{i=1}^n I_i \cap \mathbb{Q}[x_1, \dots, x_n] \right).$$

Using the algorithm of [169], we can compute the algebraic set on the right-hand side as follows:

1. For each  $1 \leq i \leq n$ , compute a set  $G_i$  of generators of the ideal  $I_i \cap \mathbb{Q}[x_1, \dots, x_n]$ .
2. Compute a zero-dimensional parametrization  $\mathcal{C}$  of the system consisting of

$$\{f^A\} \cup G_1 \cup \dots \cup G_n.$$

Such computations mimic those in [169]. The complexity of this algorithm of course depends on the algebraic elimination procedure we use. For the complexity analysis in Sections 8.3.4 and 8.4.6, we employ the geometric resolution [87].

It basically consists in computing a one-dimensional parametrization of the curve defined by  $I_i$  and next computes a zero-dimensional parametrization of the finite set obtained by intersecting this curve with the hypersurface defined by  $f = 0$ .

We call *ParametricCurve* a subroutine that, taking the polynomial  $f$  and  $1 \leq i \leq n$ , computes a one-dimensional parametrization  $\mathcal{G}_i$  of the curve defined above. Also, let *IntersectCurve* be

a subroutine that, given a one-dimensional rational parametrization  $\mathcal{G}_i$  and  $f$ , outputs a zero-dimensional parametrization  $\mathcal{C}_i$  of their intersection.

Finally, we use a subroutine `Intersection` that, from the parametrizations  $\mathcal{C}_i$ 's, computes a zero-dimensional parametrization of  $\cap_{i=1}^n Z(\mathcal{C}_i)$ . The output  $\mathcal{C}$  of `Candidates` is obtained by reversing the change of variables.

---

**Algorithm 8.1:** Algorithm Candidates

---

**Input:**  $f \in \mathbb{Q}[x_1, \dots, x_n]$ ,  $A \in \text{GL}(n, \mathbb{Q})$   
**Output:** A zero-dimensional parametrization  $\mathcal{C}$

- 1 **for**  $1 \leq i \leq n$  **do**
- 2      $\mathcal{G}_i \leftarrow \text{ParametricCurve}(f^A, i)$
- 3      $\mathcal{C}_i \leftarrow \text{IntersectCurve}(\mathcal{G}_i, f)$
- 4  $\mathcal{C} \leftarrow \text{Intersection}(\mathcal{C}_1, \dots, \mathcal{C}_n)$
- 5  $\mathcal{C} \leftarrow \mathcal{C}^{-A}$
- 6 **return**  $\mathcal{C}$

---

## 8.3 The algorithm using roadmaps

### 8.3.1 Simplification

We introduce in this subsection a method to reduce our problem to the case where  $\mathcal{H} \cap \mathbb{R}^n$  is bounded for all  $\mathbf{x} \in \mathbb{R}^n$ . Such assumptions are required to prove the results in Subsection 8.3.2. Our technique is inspired by [9, Section 12.6]. The idea is to associate to the possibly unbounded algebraic set  $\mathcal{H} \cap \mathbb{R}^n$  a bounded real algebraic set whose isolated points are strongly related to  $\mathcal{I}(\mathcal{H})$ . The construction of such an algebraic set is as follows.

Let  $x_{n+1}$  be a new variable and  $0 < \rho \in \mathbb{R}$  such that  $\rho$  is greater than the Euclidean norm  $\|\cdot\|$  of every isolated point of  $\mathcal{H} \cap \mathbb{R}^n$ . Note that such a  $\rho$  can be obtained from a finite set of points containing the isolated points of  $\mathcal{H} \cap \mathbb{R}^n$ . We explain in Subsection 8.2.2 how to compute such a finite set.

We consider the algebraic set  $\mathcal{V}$  defined by the system

$$f = 0, \quad x_1^2 + \dots + x_n^2 + x_{n+1}^2 - \rho^2 = 0.$$

Let  $\psi$  be the projection  $(x_1, \dots, x_n, x_{n+1}) \mapsto (x_1, \dots, x_n)$ . The real counterpart of  $\mathcal{V}$  is the intersection of  $\mathcal{H}$  lifted to  $\mathbb{R}^{n+1}$  with the sphere of center  $\mathbf{0}$  and radius  $\rho$ . Therefore,  $\mathcal{V}$  is a bounded real algebraic set in  $\mathbb{R}^{n+1}$ . Moreover, the restriction of  $\psi$  to  $\mathcal{V} \cap \mathbb{R}^{n+1}$  is exactly  $\mathcal{H} \cap B(\mathbf{0}, \rho)$ . By the definition of  $\rho$ , this image contains all the real isolated points of  $\mathcal{H}$ . Lemma 8.3.1 below relates  $\mathcal{I}(\mathcal{H})$  to the isolated points of  $\mathcal{V} \cap \mathbb{R}^{n+1}$ .

**Lemma 8.3.1.** *Let  $\mathcal{V}$  and  $\psi$  as above. We denote by  $\mathcal{I}(\mathcal{V}) \subset \mathbb{R}^{n+1}$  the set of real isolated points of  $\mathcal{V}$  with non-zero  $x_{n+1}$  coordinate. Then,  $\psi(\mathcal{I}(\mathcal{V})) = \mathcal{I}(\mathcal{H})$ .*

*Proof.* Note that  $\psi(\mathcal{V} \cap \mathbb{R}^{n+1}) = (\mathcal{H} \cap \mathbb{R}^n) \cap B(\mathbf{0}, \rho)$ . We consider a real isolated point  $\mathbf{x}' = (\alpha_1, \dots, \alpha_n, \alpha_{n+1})$  of  $\mathcal{V}$  with  $\alpha_{n+1} \neq 0$  and  $\mathbf{x} = \psi(\mathbf{x}') = (\alpha_1, \dots, \alpha_n)$ . Assume by contradiction that  $\mathbf{x} \notin \mathcal{I}(\mathcal{H})$ , we will prove that  $\mathbf{x}' \notin \mathcal{I}(\mathcal{V})$ , i.e., for any  $r > 0$ , there exists  $\mathbf{y}' = (\beta_1, \dots, \beta_n, \beta_{n+1}) \in \mathcal{V} \cap \mathbb{R}^{n+1}$  such that  $\|\mathbf{y}' - \mathbf{x}'\| < r$ . Since  $\mathbf{x}$  is not isolated, there exists a point  $\mathbf{y} \neq \mathbf{x}$  such that

$$\|\mathbf{y} - \mathbf{x}\| < \frac{r}{1 + 2\rho/|\alpha_{n+1}|}.$$

Let  $\mathbf{y}' \in \psi^{-1}(\mathbf{y})$  such that  $\alpha_{n+1}\beta_{n+1} \geq 0$ . We have that  $\|\mathbf{x}\|^2 + \alpha_{n+1}^2 = \|\mathbf{y}\|^2 + \beta_{n+1}^2 = \rho^2$ . Now we estimate

$$\begin{aligned} \|\mathbf{y}\|^2 - \|\mathbf{x}\|^2 &= (\|\mathbf{x}\| + \|\mathbf{y}\|) \cdot \|\mathbf{y}\| - \|\mathbf{x}\| \leq 2\rho \cdot \|\mathbf{y} - \mathbf{x}\|, \\ |\alpha_{n+1} - \beta_{n+1}| &\leq \frac{|\alpha_{n+1}^2 - \beta_{n+1}^2|}{|\alpha_{n+1}|} = \frac{\|\mathbf{y}\|^2 - \|\mathbf{x}\|^2}{|\alpha_{n+1}|} \leq \frac{2\rho \cdot \|\mathbf{y} - \mathbf{x}\|}{|\alpha_{n+1}|}. \end{aligned}$$

Finally,

$$\|\mathbf{y}' - \mathbf{x}'\| \leq \|\mathbf{y} - \mathbf{x}\| + |\alpha_{n+1} - \beta_{n+1}| \leq \left(1 + \frac{2\rho}{|\alpha_{n+1}|}\right) \|\mathbf{y} - \mathbf{x}\| < r.$$

So,  $\mathbf{x}'$  is not isolated in  $\mathcal{V} \cap \mathbb{R}^{n+1}$ . This contradiction implies that  $\psi(\mathcal{I}(\mathcal{V})) \subset \mathcal{I}(\mathcal{H})$ .

It remains to prove that  $\mathcal{I}(\mathcal{H}) \subset \psi(\mathcal{I}(\mathcal{V}))$ . For any  $\mathbf{x} \in \mathcal{I}(\mathcal{H})$ , we consider a ball

$$B(\mathbf{x}, r') \subset B(\mathbf{0}, \rho) \subset \mathbb{R}^n$$

such that  $B(\mathbf{x}, r') \cap \mathcal{H} = \{\mathbf{x}\}$ . We have that

$$\psi^{-1}(B(\mathbf{x}, r')) \cap \mathcal{V} \cap \mathbb{R}^{n+1} = \psi^{-1}(\mathbf{x}) \cap \mathcal{V} \cap \mathbb{R}^{n+1},$$

which is finite. So, all the points in  $\psi^{-1}(B(\mathbf{x}, r')) \cap \mathcal{V} \cap \mathbb{R}^{n+1}$  are isolated. Since  $\mathcal{I}(\mathcal{H}) \subset B(\mathbf{0}, \rho)$ , we deduce that  $\mathcal{I}(\mathcal{H})$  is contained in  $\psi(\mathcal{I}(\mathcal{V}))$ .

Thus, we conclude that  $\psi(\mathcal{I}(\mathcal{V})) = \mathcal{I}(\mathcal{H})$ .  $\square$

Note that the condition  $x_{n+1} \neq 0$  is crucial. For a connected component  $\mathcal{C}$  of  $\mathcal{H} \cap \mathbb{R}^n$  that is not a singleton, its intersection with the closed ball  $\overline{B(\mathbf{0}, \rho)}$  can have an isolated point on the boundary of the ball, which corresponds to an isolated point of  $\mathcal{V} \cap \mathbb{R}^{n+1}$ . This situation depends on the choice of  $\rho$  and can be easily detected by checking the vanishing of the coordinate  $x_{n+1}$ .

### 8.3.2 Geometric results

By Proposition 8.2.4, the real points of

$$\bigcap_{i=1}^n \lim_{\varepsilon} \text{crit}(\pi_i, \mathcal{H}_\varepsilon)$$

are potential isolated points of  $\mathcal{H} \cap \mathbb{R}^n$ . We study now how to identify, among those candidates, which points are truly isolated.

We use the same  $g = x_1^2 + \dots + x_{n+1}^2 - \rho^2$  and  $\mathcal{V} = V(f, g) \subset \mathbb{C}^{n+1}$  as in Subsection 8.3.1. Let  $\mathcal{V}_\varepsilon = V(f^2 - \varepsilon^2, g) \subset \mathbb{C}\langle\varepsilon\rangle^{n+1}$ ; it is the union of two algebraic sets defined respectively by  $f - \varepsilon = g = 0$  and  $f + \varepsilon = g = 0$ .

**Lemma 8.3.2.** *Let  $\mathbf{x} \in \mathcal{V} \cap \mathbb{R}^{n+1}$  such that its  $x_{n+1}$ -coordinate is non-zero. Then,  $\mathbf{x}$  is not an isolated point of  $\mathcal{V} \cap \mathbb{R}^{n+1}$  if and only if there exists a semi-algebraically connected component  $\mathcal{C}_\varepsilon$  of  $\mathcal{V}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^{n+1}$ , bounded over  $\mathbb{R}$ , such that  $\{\mathbf{x}\} \subsetneq \lim_\varepsilon \mathcal{C}_\varepsilon$ .*

*Proof.* Let  $\mathbf{x} = (\alpha_1, \dots, \alpha_{n+1}) \in \mathcal{V} \cap \mathbb{R}^{n+1}$  such that  $\alpha_{n+1} \neq 0$ . As  $f(\alpha_1, \dots, \alpha_n) = 0$ , by Lemma 8.2.3, there exists a point  $\mathbf{x}_\varepsilon = (\beta_1, \dots, \beta_{n+1}) \in \mathbb{R}\langle\varepsilon\rangle^{n+1}$  such that  $(\beta_1, \dots, \beta_n) \in \mathcal{H}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^n$  and  $\lim_\varepsilon (\beta_1, \dots, \beta_n) = (\alpha_1, \dots, \alpha_n)$ . Since  $\alpha_{n+1} \neq 0$ , we can choose  $\beta_{n+1}$  such that  $g(\mathbf{x}_\varepsilon) = 0$ . Therefore, for any  $\mathbf{x}$  as above, there exists  $\mathbf{x}_\varepsilon \in \mathcal{V}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^{n+1}$  such that  $\lim_\varepsilon \mathbf{x}_\varepsilon = \mathbf{x}$ .

Since  $\mathcal{V}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^{n+1}$  lies on the sphere (in  $\mathbb{R}\langle\varepsilon\rangle^{n+1}$ ) defined by  $g = 0$ , every connected component of  $\mathcal{V}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^{n+1}$  is bounded over  $\mathbb{R}$ . Hence, the points of  $\mathcal{V} \cap \mathbb{R}^{n+1}$  whose  $x_{n+1}$ -coordinates are not zero are contained in  $\lim_\varepsilon \mathcal{V}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^{n+1}$ .

Let  $\mathbf{x}$  be a non-isolated point of  $\mathcal{V} \cap \mathbb{R}^{n+1}$  whose  $x_{n+1}$ -coordinate is not zero. We assume by contradiction that for any semi-algebraically connected component  $\mathcal{C}_\varepsilon$  of  $\mathcal{V}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^{n+1}$  (which is bounded over  $\mathbb{R}$  by above), then it happens that either  $\lim_\varepsilon \mathcal{C}_\varepsilon = \{\mathbf{x}\}$  or  $\mathbf{x} \notin \lim_\varepsilon \mathcal{C}_\varepsilon$ .

Since  $\mathcal{V}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^{n+1}$  has finitely many connected components, the number of connected components of the second type is also finite. Since  $\mathcal{V} \cap \mathbb{R}^{n+1}$  is not a singleton (by the existence of  $\mathbf{x}$ ), the connected components of the second type exist. So, we enumerate them as  $\mathcal{C}_1, \dots, \mathcal{C}_k$  and  $\mathbf{x} \notin \lim_\varepsilon \mathcal{C}_j$  for  $1 \leq j \leq k$ .

As  $\mathbf{x}$  is not isolated in  $\mathcal{V} \cap \mathbb{R}^{n+1}$  with non-zero  $x_{n+1}$ -coordinate by assumption, there exists a sequence of points  $(\mathbf{x}_i)_{i \geq 0}$  in  $\mathcal{V} \cap \mathbb{R}^{n+1}$  of non-zero  $x_{n+1}$ -coordinates that converges to  $\mathbf{x}$ . Since there are finitely many  $\mathcal{C}_i$ , there exists an index  $j$  such that  $\lim_\varepsilon \mathcal{C}_j$  contains a subsequence of  $(\mathbf{x}_i)_{i \geq 0}$ . By [9, Proposition 12.49], the limit of the semi-algebraically connected component  $\mathcal{C}_j$  (which is bounded over  $\mathbb{R}$ ) is a closed and connected semi-algebraic set. It follows that  $\mathbf{x} \in \lim_\varepsilon \mathcal{C}_j$ , which is a contradiction. Therefore, there exists a semi-algebraically connected component of  $\mathcal{V}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^{n+1}$ , bounded over  $\mathbb{R}$ , such that  $\{\mathbf{x}\} \subsetneq \lim_\varepsilon \mathcal{C}_\varepsilon$ .

It remains to prove the reverse implication. Assume that  $\{\mathbf{x}\} \subsetneq \lim_\varepsilon \mathcal{C}_\varepsilon$  for some semi-algebraically connected component  $\mathcal{C}_\varepsilon$  of  $\mathcal{V}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^{n+1}$  that is bounded over  $\mathbb{R}$ . As  $\lim_\varepsilon \mathcal{C}_\varepsilon$  is connected, we finish the proof.  $\square$

**Lemma 8.3.3.** *Let  $\mathbf{x} \in \mathcal{V} \cap \mathbb{R}^{n+1}$  whose  $x_{n+1}$ -coordinate is non-zero. Assume that  $\mathbf{x}$  is not an isolated point of  $\mathcal{V} \cap \mathbb{R}^{n+1}$ . For any semi-algebraically connected component  $\mathcal{C}_\varepsilon$  of  $\mathcal{V}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^{n+1}$ , bounded over  $\mathbb{R}$ , such that  $\{\mathbf{x}\} \subsetneq \lim_\varepsilon \mathcal{C}_\varepsilon$ , there exists  $1 \leq i \leq n$  such that  $\mathcal{C}_\varepsilon \cap \text{crit}(\pi_i, \mathcal{V}_\varepsilon)$  contains a point  $\mathbf{x}'_\varepsilon$  which satisfies  $\lim_\varepsilon \mathbf{x}'_\varepsilon \neq \mathbf{x}$ .*

*Proof.* Let  $\mathcal{C}_\varepsilon$  be semi-algebraically connected component of  $\mathcal{V}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^{n+1}$ , bounded over  $\mathbb{R}$ , such that  $\{\mathbf{x}\} \subsetneq \lim_\varepsilon \mathcal{C}_\varepsilon$ . Lemma 8.3.2 ensures the existence of such a connected component  $\mathcal{C}_\varepsilon$ .

Now let  $\mathbf{x}_\varepsilon$  and  $\mathbf{y}_\varepsilon$  be two points contained in  $\mathcal{C}_\varepsilon$  such that  $\lim_\varepsilon \mathbf{x}_\varepsilon = \mathbf{x}$ ,  $\lim_\varepsilon \mathbf{y}_\varepsilon = \mathbf{y}$  and  $\mathbf{x} \neq \mathbf{y}$ . Let  $\mathbf{x} = (\alpha_1, \dots, \alpha_{n+1})$  and  $\mathbf{y} = (\beta_1, \dots, \beta_{n+1})$ . Since  $\mathbf{x} \neq \mathbf{y}$ , there exists  $1 \leq i \leq n+1$  such that  $\alpha_i \neq \beta_i$ . Note that if  $(\alpha_1, \dots, \alpha_n) = (\beta_1, \dots, \beta_n)$  for any  $\mathbf{y} \in \lim_\varepsilon \mathcal{C}_\varepsilon$ , then  $\lim_\varepsilon \mathcal{C}_\varepsilon$  contains at most two points (by the constraint  $g = 0$ ). However, since  $\lim_\varepsilon \mathcal{C}_\varepsilon$  is connected and contains at least two points, it must be an infinite set. So, we can choose  $\mathbf{y}$  such that  $(\alpha_1, \dots, \alpha_n) \neq (\beta_1, \dots, \beta_n)$ .

As  $\mathcal{C}_\varepsilon$  is closed in  $\mathbb{R}\langle\varepsilon\rangle^{n+1}$  (as a connected component of an algebraic set) and bounded over  $\mathbb{R}$  by definition, its projection on the  $x_i$ -coordinate is a closed interval  $[a, b] \subset \mathbb{R}\langle\varepsilon\rangle$  (see [9, Theorem 3.23]), which is bounded over  $\mathbb{R}$  (because  $\mathcal{C}_\varepsilon$  is). Also, since  $[a, b]$  is closed, there exist  $\mathbf{x}'_a$  and  $\mathbf{x}'_b$  in  $\mathbb{R}\langle\varepsilon\rangle^{n+1}$  such that  $\mathbf{x}'_a \in \pi_i^{-1}(a) \cap \mathcal{C}_\varepsilon \cap \text{crit}(\pi_i, \mathcal{V}_\varepsilon)$  and  $\mathbf{x}'_b \in \pi_i^{-1}(b) \cap \mathcal{C}_\varepsilon \cap \text{crit}(\pi_i, \mathcal{V}_\varepsilon)$ . Since  $\alpha_i \neq \beta_i$  both lying in  $\mathbb{R}$ ,  $\{\alpha_i, \beta_i\} \subset [\lim_\varepsilon a, \lim_\varepsilon b]$  implies that  $\lim_\varepsilon a \neq \lim_\varepsilon b$ . It follows that  $\lim_\varepsilon \mathbf{x}'_a \neq \lim_\varepsilon \mathbf{x}'_b$ . Thus, at least one point among  $\lim_\varepsilon \mathbf{x}'_a$  and  $\lim_\varepsilon \mathbf{x}'_b$  does not coincide with  $\mathbf{x}$ . Hence, there exists a point  $\mathbf{x}'_\varepsilon$  in  $\mathcal{C}_\varepsilon \cap \text{crit}(\pi_i, \mathcal{V}_\varepsilon)$  such that  $\lim_\varepsilon \mathbf{x}'_\varepsilon \neq \mathbf{x}$ .  $\square$

We can easily deduce from Lemma 8.3.2 and Lemma 8.3.3 the following proposition, which is the main ingredient of our algorithm.

**Proposition 8.3.4.** *Let  $\mathbf{x} \in \bigcap_{i=1}^n \lim_\varepsilon \text{crit}(\pi_i, \mathcal{V}_\varepsilon) \cap \mathbb{R}\langle\varepsilon\rangle_b^{n+1}$  whose  $x_{n+1}$ -coordinate is non-zero. Then,  $\mathbf{x}$  is not an isolated point of  $\mathcal{V} \cap \mathbb{R}^{n+1}$  if and only if there exist  $1 \leq i \leq n$  and a connected component  $\mathcal{C}_\varepsilon$  of  $\mathcal{V}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^{n+1}$ , which is bounded over  $\mathbb{R}$ , such that  $\mathcal{C}_\varepsilon \cap \text{crit}(\pi_i, \mathcal{H}_\varepsilon)$  contains  $\mathbf{x}_\varepsilon$ ,  $\mathbf{x}'_\varepsilon$  satisfying  $\mathbf{x} = \lim_\varepsilon \mathbf{x}_\varepsilon \neq \lim_\varepsilon \mathbf{x}'_\varepsilon$ .*

### 8.3.3 Description of the algorithm

Our algorithm takes as input a polynomial  $f \in \mathbb{Q}[x_1, \dots, x_n]$  and computes a data consisting of:

- A zero-dimensional parametrization  $\mathcal{C}$  such that  $\mathcal{I}(\mathcal{H})$  is a subset of the zeros of  $\mathcal{C}$ .
- A finite set  $\mathcal{B}$  of isolating intervals of solutions of  $w(t)$  that corresponds to  $\mathcal{I}(\mathcal{H})$ .

The first step computes a zero-dimensional parametrization  $\mathcal{C}$  encoding a finite set of points which contains  $\mathcal{I}(\mathcal{H})$ . This is done by the subroutine Candidates described in Subsection 8.2.2.

The next step consists of identifying those of the candidates which are isolated in  $\mathcal{H} \cap \mathbb{R}^n$ . This step relies on Proposition 8.3.4. To reduce our problem to the context where Proposition 8.3.4 can be applied, we use Lemma 8.3.1. One needs to compute  $\rho \in \mathbb{R}$ , such that  $\rho$  is larger than the maximum norm of the real isolated points we want to compute. This value of  $\rho$  can be easily obtained by isolating the real roots of the zero-dimensional parametrization encoding the candidates. Further, we call GetNormBound a subroutine which takes as input  $\mathcal{C}$  and returns  $\rho$  as we just sketched.

We let  $g = x_1^2 + \dots + x_n^2 + x_{n+1}^2 - \rho^2$ . By Lemma 8.3.1,  $\mathcal{I}(\mathcal{H})$  is the projection of the set of real isolated points of the algebraic set  $\mathcal{V}$  defined by  $f = g = 0$  at which  $x_{n+1} \neq 0$ . Let  $\mathcal{X}$  be the set of points of  $\mathcal{V}$  projecting to the candidates encoded by  $\mathcal{C}$ .

The remain of this subsection is devoted to describe the identification of the real isolated points among the candidates  $\mathfrak{C} \cap \mathbb{R}^n$ . Let

$$\mathcal{P} = \{\mathbf{x} = (x_1, \dots, x_{n+1}) \in \mathbb{R}^{n+1} \mid (x_1, \dots, x_n) \in \mathfrak{C}, g(\mathbf{x}) = 0, x_{n+1} \neq 0\}.$$

Proposition 8.3.4 would lead us to compute  $\text{crit}(\pi_i, \mathcal{V}_\varepsilon^A)$  as well as a roadmap of  $\mathcal{V}_\varepsilon^A$ . As explained in the introduction, this induces computations over the ground field  $\mathbb{R}\langle\varepsilon\rangle$  which we want to avoid. Hence, in what follows, we replace the infinitesimal  $\varepsilon$  by a sufficiently small  $e \in \mathbb{R}$  then adapt the results of Subsection 8.3.2 to  $\mathcal{V}_e$ .

Let  $t$  be a new variable,

$$\mathcal{V}_t = \{(\mathbf{x}, t) \in \mathbb{R}^{n+1} \times \mathbb{R} \mid f^A(\mathbf{x}) = t, g(\mathbf{x}) = 0\},$$

$\pi_{\mathbf{x}} : (\mathbf{x}, t) \mapsto \mathbf{x}$  and  $\pi_t : (\mathbf{x}, t) \mapsto t$ . Note that  $\mathcal{V}_t$  is smooth. Recall that the set of critical values of the restriction of  $\pi_t$  to  $\mathcal{V}_t$  is finite by the algebraic Sard's theorem (Theorem 2.5.11). Given a semi-algebraic set  $\mathcal{S} \subset \mathbb{R}^{n+1} \times \mathbb{R}$  in the  $(\mathbf{x}, t)$ -coordinates and a subset  $\mathcal{I}$  of  $\mathbb{R}$ , the notation  $\mathcal{S}_{\mathcal{I}}$  stands for the fiber  $\pi_t^{-1}(\mathcal{I}) \cap \mathcal{S}$ .

Let  $\mathcal{V}_e \in \mathbb{C}^{n+1}$  denote the algebraic set defined by

$$(f^A - e) \cdot (f^A + e) = g = 0$$

for a number  $e \in \mathbb{R}$ . By definition,  $\mathcal{V}_e \cap \mathbb{R}^{n+1}$  is compact for any  $e \in \mathbb{R}$ . Hence, the restriction of  $\pi_t$  to  $\mathcal{V}_t$  is proper. Then, by Thom's isotopy lemma [47],  $\pi_t$  realizes a locally trivial fibration over any open connected subset of  $\mathbb{R}$  which does not intersect the set of critical values of the restriction of  $\pi_t$  to  $\mathcal{V}_t$ . Let  $\eta \in \mathbb{R}$  such that the open set  $] - \eta, 0[ \cup ] 0, \eta[$  does not contain any critical value of the restriction of  $\pi_t$  to the algebraic set  $\mathcal{V}_t$ . Hence,  $\mathcal{V}_e$  is non-singular for  $e \in ] 0, \eta[$ ,  $(\mathcal{V}_e \cap \mathbb{R}^{n+1}) \times (] - \eta, 0[ \cup ] 0, \eta[)$  is diffeomorphic to  $\mathcal{V}_{t \in ] - \eta, 0[ \cup ] 0, \eta[}$ .

We need to mention that  $\text{crit}(\pi_i, \mathcal{H}_e)$  corresponds to the critical points of  $\pi_i$  restricted to  $\mathcal{V}_e$  with non-zero  $x_{n+1}$ -coordinate. Further, we use  $\text{crit}(\pi_i, \mathcal{V}_e)$  to address those latter critical points.

Now, for  $1 \leq i \leq n$ , we define  $\mathcal{W}_i$  as the closure of

$$\mathcal{V}_t \cap \left\{ (\mathbf{x}, t) \in \mathbb{R}^{n+1} \times \mathbb{R} \mid \frac{\partial f}{\partial x_i}(\mathbf{x}) \neq 0, \frac{\partial f}{\partial x_j}(\mathbf{x}) = 0 \text{ for } j \neq i, x_{n+1} \neq 0 \right\}.$$

Since  $A$  is assumed to be generically chosen,  $\mathcal{W}_i$  is either empty or one-equidimensional (because

$$\left\langle y \cdot \frac{\partial f}{\partial x_i} - 1, \frac{\partial f}{\partial x_j} \text{ for } j \neq i \right\rangle$$

either defines an empty set or a one-equidimensional algebraic set by [169]). This implies that the set of singular points of  $\mathcal{W}_i$  is finite.

By [116], the set of non-properness of the restriction of  $\pi_t$  to  $\mathcal{W}_i$  is finite. Using again [116], the restriction of  $\pi_t$  to  $\mathcal{W}_i$  realizes a locally trivial fibration over any connected open subset which does not meet the union of the images by  $\pi_t$  of the singular points of  $\mathcal{W}_i$ , the set of non-properness,

and the set of critical values of the restriction of  $\pi_t$  to  $\mathcal{W}_i$ . We let  $\eta'_i$  be the minimum of the absolute values of the points in this union.

We choose now  $0 < e_0 < \min(\eta, \eta'_1, \dots, \eta'_n)$ . We call `SpecializationValue` a subroutine that takes as input  $f$  and  $g$  and returns such a rational number  $e_0$ . Note that `SpecializationValue` is easily obtained from elimination algorithms solving polynomial systems (from which we can compute critical values) and from [172] to compute the set of non-properness of some map.

With  $e_0$  as above, we denote  $\mathcal{I} = ]-e_0, 0[ \cup ]0, e_0[$ . Let  $\mathcal{W}_{i,\mathcal{I}}$  is semi-algebraically diffeomorphic to  $\mathcal{W}_{i,e} \times \mathcal{I}$  for every  $e \in \mathcal{I}$ . As  $\mathcal{V}_e$  is nonsingular, the critical locus  $\text{crit}(\pi_i, \mathcal{V}_e)$  is guaranteed to be finite by the genericity of the change of variables  $A$  (hence  $\mathcal{W}_{i,e}$  is) and that  $\text{crit}(\pi_i, \mathcal{V}_e) \cap \mathbb{R}^{n+1}$  coincides with  $\pi_{\mathbf{x}}(\mathcal{W}_{i,e})$ . Thus, the above diffeomorphism implies that, for any connected component  $\mathcal{C}$  of  $\mathcal{W}_{i,\mathcal{I}}$ ,  $\mathcal{C}$  is diffeomorphic to an open interval in  $\mathbb{R}$ . Moreover, if  $\mathcal{C}$  is bounded, then  $\overline{\mathcal{C}} \setminus \mathcal{C}$  contains exactly two points which satisfy respectively  $f = 0$  and  $f^2 = e_0^2$ . We now consider

$$\mathcal{L}_i = \left\{ \mathbf{x} \in \mathbb{R}^{n+1} \mid 0 < f < e_0, g = 0, \frac{\partial f}{\partial x_j} = 0 \text{ for } j \neq i, x_{n+1} \neq 0 \right\}.$$

It is the intersection of the Zariski closure  $\mathcal{K}_i$  of the solution set to

$$\left\{ \frac{\partial f}{\partial x_i} \neq 0, \frac{\partial f}{\partial x_j} = 0 \text{ for } j \neq i, x_{n+1} \neq 0 \right\}$$

with the semi-algebraic set defined by  $0 < f < e_0$ . Note that  $\mathcal{K}_i$  is either empty or one-equidimensional. As  $\mathcal{V}_e$  is nonsingular for  $e \in \mathcal{I}$ ,  $\mathcal{L}_i$  and  $\mathcal{L}_j$  are disjoint for  $i \neq j$ . Since the restriction of  $\pi_{\mathbf{x}}$  to  $\mathcal{V}_t$  is an isomorphism between the algebraic sets  $\mathcal{V}_t$  and  $\mathbb{R}^{n+1}$  with the inverse map  $\mathbf{x} \mapsto (\mathbf{x}, f(\mathbf{x}))$ , the properties of  $\mathcal{W}_{\mathcal{I}}$  mentioned above are transferred to its image  $\mathcal{L}_i$  by the projection  $\pi_{\mathbf{x}}$ .

Further, we consider a subroutine `ParametricCurve` which takes as input  $f$  and  $i \in [1, n]$  and returns a rational parametrization  $\mathfrak{K}_i$  of  $\mathcal{K}_i$ . Also, let `Union` be a subroutine that takes a family of rational parametrizations  $\mathfrak{K}_1, \dots, \mathfrak{K}_n$  to compute a rational parametrization encoding the union of the algebraic curves defined by the  $\mathfrak{K}_i$ 's. We denote by  $\mathfrak{K}$  the output of `Union`; it encodes  $\mathcal{K} = \cup_{i=1}^n \mathcal{K}_i$ . We refer to [174, Appendix J.2] for these two subroutines.

Lemma 8.3.5 below establishes a *well-defined* notion of limit for a point  $\mathbf{x}_e \in \text{crit}(\pi_i, \mathcal{V}_e)$  when  $e$  tends to 0.

**Lemma 8.3.5.** *Let  $e_0$  and  $\mathcal{L}_i$  be as above. For  $e \in ]0, e_0[$  and  $\mathbf{x}_e \in \text{crit}(\pi_i, \mathcal{V}_e) \cap \mathbb{R}^{n+1}$ , there exists a (unique) connected component  $\mathcal{C}$  of  $\mathcal{L}_i$  containing  $\mathbf{x}_e$ . If  $\mathcal{C}$  is bounded, let  $\mathbf{x}$  be the only point in  $\overline{\mathcal{C}}$  satisfying  $f(\mathbf{x}) = 0$ , then  $\mathbf{x} \in \lim_{\varepsilon} \text{crit}(\pi_i, \mathcal{V}_{\varepsilon}) \cap \mathbb{R}\langle \varepsilon \rangle^{n+1}$ . Thus, we set  $\lim_0 \mathbf{x}_e = \mathbf{x}$ .*

*Moreover, the extension  $\text{ext}(\mathcal{C}, \mathbb{R}\langle \varepsilon \rangle)$  contains exactly one point  $\mathbf{x}_{\varepsilon}$  such that  $f(\mathbf{x}_{\varepsilon})^2 = \varepsilon^2$  and  $\lim_{\varepsilon} \mathbf{x}_{\varepsilon} = \mathbf{x}$ .*

*Proof.* Since  $\mathbf{x}_e \in \text{crit}(\pi_i, \mathcal{V}_e) \cap \mathbb{R}^{n+1}$  and  $0 < e < e_0$ , we have  $\mathbf{x}_e \in \mathcal{L}_i$ , the existence of  $\mathcal{C}$  follows naturally. Let  $\mathbf{x}$  be the unique point of  $\overline{\mathcal{C}}$  satisfying  $f = 0$ . Then, the notion  $\lim_0$  is well-defined.

From the proof of [9, Theorem 12.43], we have that

$$\lim_{\varepsilon} \text{crit}(\pi_i, \mathcal{V}_{\varepsilon}) \cap \mathbb{R}\langle\varepsilon\rangle^{n+1} = \pi_{\mathbf{x}} \left( \overline{\mathcal{W}_{(0,+\infty)}} \cap V(t) \right).$$

As  $\pi_{\mathbf{x}} \left( \overline{\mathcal{W}_{(0,+\infty)}} \cap V(t) \right)$  is the set of points corresponding to  $f = 0$  of  $\mathcal{L}_i$ , we deduce that  $\mathbf{x} \in \lim_{\varepsilon} \text{crit}(\pi_i, \mathcal{V}_{\varepsilon}) \cap \mathbb{R}\langle\varepsilon\rangle^{n+1}$ .

Since the extension  $\text{ext}(\mathcal{C}, \mathbb{R}\langle\varepsilon\rangle)$  is a connected component of  $\text{ext}(\mathcal{L}_i, \mathbb{R}\langle\varepsilon\rangle)$  and homeomorphic to an open interval in  $\mathbb{R}\langle\varepsilon\rangle$ , there exists  $\mathbf{x}_{\varepsilon} \in \text{ext}(\mathcal{C}, \mathbb{R}\langle\varepsilon\rangle)$  such that  $f(\mathbf{x}_{\varepsilon})^2 = \varepsilon^2$ . Moreover, since  $0 = \lim_{\varepsilon} f(\mathbf{x}_{\varepsilon})^2 = f(\lim_{\varepsilon} \mathbf{x}_{\varepsilon})^2$  and  $\mathbf{x}$  is the only point in  $\overline{\mathcal{C}}$  satisfying  $f = 0$ , we conclude that  $\lim_{\varepsilon} \mathbf{x}_{\varepsilon} = \mathbf{x}$ .  $\square$

Now, let  $\mathcal{R}_e$  be a roadmap associated to the algebraic set  $\mathcal{V}_e$ , i.e.  $\mathcal{R}_e$  is contained in  $\mathcal{V}_e \cap \mathbb{R}^{n+1}$ , of at most dimension one and has non-empty intersection with every connected component of  $(\mathcal{V}_e \cup \mathcal{V}_{-e}) \cap \mathbb{R}^{n+1}$ . We also require that  $\mathcal{R}_e$  contains  $\cup_{i=1}^n (\text{crit}(\pi_i, \mathcal{V}_e) \cup \text{crit}(\pi_i, \mathcal{V}_{-e})) \cap \mathbb{R}^{n+1}$ . The proposition below is the key of our algorithm.

**Proposition 8.3.6.** *Given  $e \in ]0, e_0[$  and  $\mathcal{I} = ] - e_0, 0[ \cup ]0, e_0[$  as above. Let  $\mathcal{L} = \cup_{i=1}^n \mathcal{L}_i$  and  $\mathbf{x} \in \mathcal{P}$ . Then  $\mathbf{x}$  is not isolated in  $\mathcal{V} \cap \mathbb{R}^{n+1}$  if and only if there exists  $\mathbf{x}' \in \mathcal{P}$  such that  $\mathbf{x}$  and  $\mathbf{x}'$  are connected in  $\mathcal{P} \cup \mathcal{L} \cup \mathcal{R}_e$ .*

*Proof.* Assume first that  $\mathbf{x}$  is not isolated. By Proposition 8.3.4, there exists  $1 \leq i \leq n$  and a connected component  $\mathcal{C}_{\varepsilon}$  of  $\mathcal{V}_{\varepsilon} \cap \mathbb{R}\langle\varepsilon\rangle^{n+1}$ , which is bounded over  $\mathbb{R}$ , such that  $\mathcal{C}_{\varepsilon} \cap \text{crit}(\pi_i, \mathcal{V}_{\varepsilon})$  contains  $\mathbf{x}_{\varepsilon}$  and  $\mathbf{x}'_{\varepsilon}$  satisfying  $\mathbf{x} = \lim_{\varepsilon} \mathbf{x}_{\varepsilon} \neq \lim_{\varepsilon} \mathbf{x}'_{\varepsilon}$ . By the choice of  $e_0$ , there exist a diffeomorphism  $\theta : \mathcal{V}_{t, \mathcal{I}} \rightarrow \mathcal{V}_e \times \mathcal{I}$  such that  $\theta(\mathcal{W}_{i, \mathcal{I}}) = \theta(\mathcal{W}_{i, e}) \times \mathcal{I}$ . Using [9, Exercise 3.2],  $\text{ext}(\theta, \mathbb{R}\langle\varepsilon\rangle)$  is a diffeomorphism between:

$$\begin{aligned} \text{ext}(\mathcal{V}_{t, \mathcal{I}}, \mathbb{R}\langle\varepsilon\rangle) &\cong \text{ext}(\mathcal{V}_e, \mathbb{R}\langle\varepsilon\rangle) \times \text{ext}(\mathcal{I}, \mathbb{R}\langle\varepsilon\rangle), \\ \text{ext}(\mathcal{W}_{i, \mathcal{I}}, \mathbb{R}\langle\varepsilon\rangle) &\cong \text{ext}(\mathcal{W}_{i, e}, \mathbb{R}\langle\varepsilon\rangle) \times \text{ext}(\mathcal{I}, \mathbb{R}\langle\varepsilon\rangle). \end{aligned}$$

As  $\pi_{\mathbf{x}}$  is an isomorphism from  $\mathcal{V}_t$  to  $\mathbb{R}^{n+1}$ , there exists a (unique) bounded connected component  $\mathcal{C}_e$  of  $\mathcal{V}_e \cap \mathbb{R}^{n+1}$  s.t.  $\mathcal{C}_e$  is diffeomorphic to  $\text{ext}(\mathcal{C}_e, \mathbb{R}\langle\varepsilon\rangle)$ . Moreover, let  $L$  and  $L'$  be the connected components of  $\text{ext}(\mathcal{L}_i, \mathbb{R}\langle\varepsilon\rangle)$  containing  $\mathbf{x}_{\varepsilon}$  and  $\mathbf{x}'_{\varepsilon}$  respectively and  $\mathbf{x}_e$  and  $\mathbf{x}'_e$  ( $\in \text{ext}(\mathcal{C}_e, \mathbb{R}\langle\varepsilon\rangle)$ ) be the intersections of  $\text{ext}(\mathcal{C}_e, \mathbb{R}\langle\varepsilon\rangle)$  with  $L$  and  $L'$  respectively. Then,  $\lim_{\varepsilon} \mathbf{x}_{\varepsilon}$  ( $\lim_{\varepsilon} L'$ ) connects  $\lim_{\varepsilon} \mathbf{x}_e$  ( $\lim_{\varepsilon} \mathbf{x}'_e$ ) to  $\mathbf{x}$  ( $\mathbf{x}'$ ). As  $\lim_{\varepsilon} \mathbf{x}_e$  and  $\lim_{\varepsilon} \mathbf{x}'_e$  are connected in  $\mathcal{C}_e$ , we conclude that  $\mathbf{x}$  and  $\mathbf{x}'$  are also connected in  $\mathcal{P} \cup \mathcal{L} \cup \mathcal{R}_e$ . The reverse implication is immediate using the above techniques  $\square$

From Lemma 8.3.5 and Proposition 8.3.6, any  $e$  lying in the interval  $]0, e_0[$  defined above can be used to replace the infinitesimal  $\varepsilon$ . So, we simply take  $e = e_0$ . For  $1 \leq i \leq n$ , we use a subroutine ZeroDimSolve which takes as input  $\left\{ f - e_0, g, \frac{\partial f}{\partial x_j} \text{ for all } j \neq i \right\}$  to compute a zero-dimensional parametrization  $\Omega_i$  such that  $\text{crit}(\pi_i, \mathcal{V}_e) = \{ \mathbf{x} \in Z(\Omega_i) \mid x_{n+1} \neq 0 \}$ .

To use Proposition 8.3.6, we need to compute  $\mathcal{R}_{e_0/2}$ , which we refer to the algorithm Roadmap in [174]. This algorithm allows us to compute roadmaps for smooth and bounded real algebraic sets, which is indeed the case of  $\mathcal{V}_{e_0} \cap \mathbb{R}^{n+1}$ . First, we call (another) Union that, on the zero-dimensional parametrizations  $\mathcal{Q}_i$ , computes a zero-dimensional parametrization  $\mathcal{Q}$  encoding  $\cup_{i=1}^n Z(\mathcal{Q}_i)$ . Given the polynomials  $f, g$ , the value  $e_0/2$  and the parametrization  $\mathcal{Q}$ , a combination of Union and Roadmap returns a one-dimensional parametrization  $\mathfrak{R}$  representing  $\mathcal{R}_{e_0/2}$ .

Deciding connectivity over  $\mathcal{P} \cup \mathcal{L} \cup \mathcal{R}_e$  is done as follows. We use Union to compute a rational parametrization  $\mathfrak{S}$  encoding  $\mathcal{K} \cup \mathcal{R}_e$ . Then, with input  $\mathfrak{S}, \mathcal{C}, x_{n+1} \neq 0$  and the inequalities  $0 < f < e_0$ , we use Newton Puiseux expansions and cylindrical algebraic decomposition (see [55, 180]) following [173], taking advantage of the fact that polynomials involved in rational parametrizations of algebraic curves are bivariate. We denote by ConnectivityQuery the subroutine that takes those inputs and returns  $\mathcal{C}$  and isolating boxes of the points defined by  $\mathcal{C}$  which are not connected to other points of  $\mathcal{C}$ .

Once the real isolated points of  $\mathcal{V}^A$  is computed, we remove the boxes corresponding to points at which  $x_{n+1} = 0$  and project the remaining points on the  $(x_1, \dots, x_n)$ -space to obtain the isolated points of  $\mathcal{H}^A$ . This whole step uses a subroutine which we call Remove (see [174, Appendix J]). Finally, we reverse the change of variable by applying  $A^{-1}$  to get  $\mathcal{I}(\mathcal{H})$ .

We compute a roadmap  $\mathcal{R}_e$  for  $\mathcal{V}_e^A \cap \mathbb{R}^{n+1}$  for a small  $e > 0$ . From this roadmap, we construct a semi-algebraic curve  $\mathcal{K}$  containing  $\mathcal{X}$  such that  $\mathbf{x} \in \mathcal{X}$  is isolated in  $\mathcal{V}^A \cap \mathbb{R}^{n+1}$  if and only if it is not connected to any other  $\mathbf{x}' \in \mathcal{X}$  by  $\mathcal{K}$ . We call Isolated the subroutine that takes as input  $\mathcal{C}, f^A$  and  $g$  and returns  $\mathcal{C}$  with isolating boxes  $\mathcal{B}$  of the real points of defined by  $\mathcal{C}$  which are isolated in  $\mathcal{V}^A \cap \mathbb{R}^{n+1}$ .

The pseudo-code in Algorithm 8.2 below summarizes the details of our first algorithm for computing the real isolated point of an algebraic hypersurface.

---

**Algorithm 8.2: IsolatedPoints-RoadMap**

---

**Input:** A polynomial  $f \in \mathbb{Q}[x_1, \dots, x_n]$   
**Output:** A zero-dimensional parametrization  $\mathcal{C}$  and a set  $\mathcal{B}$  of isolating intervals

- 1  $A$  chosen randomly in  $\text{GL}(n, \mathbb{Q})$
- 2  $\mathcal{C} \leftarrow \text{Candidates}(f, A)$
- 3  $\rho \leftarrow \text{GetNormBound}(\mathcal{C})$
- 4  $g \leftarrow x_1^2 + \dots + x_n^2 + x_{n+1}^2 - \rho^2$
- 5  $e_0 \leftarrow \text{SpecializationValue}(f^A, g)$
- 6 **for**  $1 \leq i \leq n$  **do**
- 7      $\Omega_i \leftarrow \text{ZeroDimSolve} \left( \left\{ f^A - e_0, g, \frac{\partial f^A}{\partial x_j} \text{ for all } j \neq i \right\} \right)$
- 8      $\mathfrak{K}_i \leftarrow \text{ParametricCurve}(f^A, i)$
- 9 **end**
- 10  $\mathfrak{K} \leftarrow \text{Union}(\mathfrak{K}_1, \dots, \mathfrak{K}_n)$
- 11  $\Omega \leftarrow \text{Union}(\Omega_1, \dots, \Omega_n)$
- 12  $\mathfrak{R} \leftarrow \text{Union}(\text{RoadMap}(f^A - e_0, g, \Omega), \text{RoadMap}(f^A + e_0, g, \Omega))$
- 13  $\mathfrak{S} \leftarrow \text{Union}(\mathfrak{K}, \mathfrak{R})$
- 14  $\mathcal{B} \leftarrow \text{ConnectivityQuery}(\mathfrak{S}, \mathcal{C}, x_{n+1} \neq 0, 0 < f^A < e_0)$
- 15  $\mathcal{C}, \mathcal{B} \leftarrow \text{Removes}(\mathcal{C}, \mathcal{B}, x_{n+1})$
- 16  $\mathcal{C}, \mathcal{B} \leftarrow \mathcal{C}^{A^{-1}}, \mathcal{B}^{A^{-1}}$
- 17 **return**  $(\mathcal{C}, \mathcal{B})$

---

### 8.3.4 Complexity analysis

This section is dedicated to the complexity analysis of Algorithm 8.2 for an input polynomial  $f \in \mathbb{Q}[x_1, \dots, x_n]$  of degree  $D$ . All complexity results are given in the number of arithmetic operations in  $\mathbb{Q}$ . Hereafter, we assume that a generic enough matrix  $A$  is found from a random choice.

We start with the subroutine Candidates. Since  $\text{crit}(\pi_i, \mathcal{H}_\varepsilon^A)$  is finite and defined by

$$(f^A - \varepsilon) \cdot (f^A + \varepsilon) = 0, \quad \frac{\partial f^A}{\partial x_j} = 0 \text{ for all } j \neq i,$$

its degree is bounded by  $2D(D-1)^{n-1}$  using Bézout bound. Consequently, the degree of the output zero-dimensional parametrization is bounded by  $2D(D-1)^{n-1}$ .

Using [169, Theorem 4] (which is based on the geometric resolution algorithm in [87]), each zero-dimensional parametrization of  $\text{crit}(\pi_i, \mathcal{H}_\varepsilon^A)$  is computed within  $O^\sim(D^{3n})$  arithmetic operations in  $\mathbb{Q}$ . The last step which takes intersections of those parametrizations is done using the algorithm in [174, Appendix J.1]; it does not change the asymptotic complexity.

We have seen that GetNormBound reduces to isolate the real roots of a zero-dimensional parametrization of degree  $D^{O(n)}$ . This can be done within  $D^{O(n)}$  operations by Uspensky's algorithm [168].

Each call to `SpecializationValue` reduces to computing critical values of  $\pi_i$  of a smooth algebraic set defined by polynomials of degree at most  $D$ . This is done using  $(nD)^{O(n)}$  arithmetic operations in  $\mathbb{Q}$  (see [91]). Using [87] for `ZeroDimSolve` and [179] for `ParametricCurve` does not increase the overall complexity. The loop is performed  $n$  times ; hence the complexity lies in  $(nD)^{O(n)}$ . All output zero-dimensional parametrizations have degree bounded by  $D^{O(n)}$ . Running `Union` on these parametrizations does not increase the asymptotic complexity. One gets then parametrizations of degree bounded by  $nD^{O(n)}$ . Finally, using [174] for `Roadmap` uses  $(nD)^{O(n \log(n))}$  arithmetic operations in  $\mathbb{Q}$  and outputs a rational parametrization of degree lying in  $(nD)^{O(n \log(n))}$ . The call to `ConnectivityQuery`, done as explained in [173] is polynomial in the degree of the roadmap.

The final steps which consist in calling `Removes` and undoing the change of variables does not change the asymptotic complexity.

Summing up altogether the above complexity estimates, one obtains an algorithm using

$$(nD)^{O(n \log(n))}$$

arithmetics operations in  $\mathbb{Q}$  at most. This ends the proof of Theorem 8.1.1.

## 8.4 The algorithm of complexity $D^{O(n)}$

This section is dedicated to design an algorithm whose arithmetic complexity lies in  $D^{O(n)}$ , thus better than the one described in the previous section. We first establish some geometric ingredients which will be used by our second algorithm.

### 8.4.1 Geometric results

Once the set of candidates  $\mathfrak{C}$  is acquired, these results allow us to identify the candidates which are the real isolated points of  $\mathcal{H}$ .

Given a candidate  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_n) \in \mathfrak{C} \cap \mathbb{R}^n$ , we want to check whether  $\boldsymbol{\eta}$  is an isolated point of  $\mathcal{H} \cap \mathbb{R}^n$ . A direct way to do so is to check whether the sphere centered at  $\boldsymbol{\eta}$  of infinitesimal radius intersects the real algebraic set  $\mathcal{H}$ . This can be done by solving the system

$$f(x_1, \dots, x_n) = \sum_{i=1}^n (x_i - \eta_i)^2 - \varepsilon = 0$$

over the field of Puiseux series  $\mathbb{R}\langle\varepsilon\rangle$ . However, this approach leads to computations involving infinitesimal, which could prevent us from obtaining a practically efficient algorithm. Therefore, we present in what follows a workaround to avoid the infinitesimals.

Let  $\mathbf{a} = (a_1, \dots, a_n)$  be a  $n$ -uple of positive rational numbers. We consider the function  $d_{\mathbf{a}} : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  defined by

$$(\mathbf{x}, \mathbf{y}) \mapsto d_{\mathbf{a}}(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^n a_i (x_i - y_i)^2}$$

where  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{y} = (y_1, \dots, y_n)$ . Note that  $d_{\mathbf{a}}$  is a metric function in  $\mathbb{R}^n$  and that it also extends to a metric over  $\mathbb{R}\langle\varepsilon\rangle^n$ . Further, we use the notations below to respectively denote spheres, open balls and closed balls with respect to the metric  $d_{\mathbf{a}}$ :

- $S(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^n \mid d_{\mathbf{a}}(\mathbf{x}, \mathbf{y}) = r\}$ ,
- $B(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^n \mid d_{\mathbf{a}}(\mathbf{x}, \mathbf{y}) < r\}$ ,
- $\bar{B}(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^n \mid d_{\mathbf{a}}(\mathbf{x}, \mathbf{y}) \leq r\}$ .

For each  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_n) \in \mathbb{C}\langle\varepsilon\rangle^n$ , we consider the map  $\delta_{\boldsymbol{\eta}}$  of distance to  $\boldsymbol{\eta}$ :

$$\begin{aligned} \delta_{\boldsymbol{\eta}} : \quad \mathbb{C}\langle\varepsilon\rangle^n &\rightarrow \mathbb{C}\langle\varepsilon\rangle, \\ \mathbf{x} = (x_1, \dots, x_n) &\mapsto \sum_{i=1}^n a_i(x_i - \eta_i)^2. \end{aligned}$$

Since  $\mathcal{H}_{\varepsilon}$  is a smooth algebraic hypersurface, by Sard's theorem [19, Theorem 9.6.2], the critical values of the restriction of  $\delta_{\boldsymbol{\eta}}$  to the algebraic set  $\mathcal{H}_{\varepsilon} \subset \mathbb{C}\langle\varepsilon\rangle^n$  is a finite subset of  $\mathbb{C}\langle\varepsilon\rangle$ . Therefore, for every  $\boldsymbol{\eta} \in \mathcal{C} \cap \mathbb{R}^n$ , the set

$$\{\delta_{\boldsymbol{\eta}}(\lim_{\varepsilon} \mathbf{x}_{\varepsilon}) \mid \mathbf{x}_{\varepsilon} \in \text{crit}(\delta_{\boldsymbol{\eta}}, \mathcal{H}_{\varepsilon}) \cap \mathbb{R}\langle\varepsilon\rangle_b^n, \lim_{\varepsilon} \mathbf{x}_{\varepsilon} \neq \boldsymbol{\eta}\}$$

is a finite set of positive elements of  $\mathbb{R}$ . Note that the above set can be empty, we use the lemma below to handle specifically this situation.

**Lemma 8.4.1.** *Assuming that there exists  $\boldsymbol{\eta} \in \mathcal{H} \cap \mathbb{R}^n$  such that the set*

$$\{\mathbf{x}_{\varepsilon} \in \text{crit}(\delta_{\boldsymbol{\eta}}, \mathcal{H}_{\varepsilon}) \cap \mathbb{R}\langle\varepsilon\rangle_b^n, \lim_{\varepsilon} \mathbf{x}_{\varepsilon} \neq \boldsymbol{\eta}\}$$

*is empty, then, exactly one among the two statements below holds:*

- i)  $\mathcal{H} \cap \mathbb{R}^n$  is connected and not bounded;
- ii)  $\mathcal{H} \cap \mathbb{R}^n$  is a single point  $\boldsymbol{\eta}$ .

*Proof.* Assume that  $\mathcal{H} \cap \mathbb{R}^n$  is the union of at least two connected components. Then, there exists a connected component of  $\mathcal{H}_{\varepsilon} \cap \mathbb{R}\langle\varepsilon\rangle^n$  that does not contain any point whose limit is  $\boldsymbol{\eta}$ . Therefore, the restriction of  $\delta_{\boldsymbol{\eta}}$  to this connected component admits a critical point over this connected component. As a consequence,  $\{\mathbf{x}_{\varepsilon} \in \text{crit}(\delta_{\boldsymbol{\eta}}, \mathcal{H}_{\varepsilon}) \cap \mathbb{R}\langle\varepsilon\rangle_b^n, \lim_{\varepsilon} \mathbf{x}_{\varepsilon} \neq \boldsymbol{\eta}\}$  is not empty, which contradicts the assumption in Lemma 8.4.1. Therefore,  $\mathcal{H} \cap \mathbb{R}^n$  has exactly one connected component.

Now we assume that  $\mathcal{H} \cap \mathbb{R}^n$  is bounded and that it is not a single point  $\boldsymbol{\eta}$ . As a consequence of Lemma 8.2.3, there exists a connected component  $\mathcal{C}_{\boldsymbol{\eta}, \varepsilon}$ , that is bounded over  $\mathbb{R}$ , of  $\mathcal{H}_{\varepsilon} \cap \mathbb{R}\langle\varepsilon\rangle^n$  such that  $\{\boldsymbol{\eta}\} \subsetneq \lim_{\varepsilon} \mathcal{C}_{\boldsymbol{\eta}, \varepsilon}$ . Thus, the distance function  $\delta_{\boldsymbol{\eta}}$  attains its maximum, which is non-zero, over  $\mathcal{H} \cap \mathbb{R}^n$ . This also contradicts the assumption that  $\{\mathbf{x}_{\varepsilon} \in \text{crit}(\delta_{\boldsymbol{\eta}}, \mathcal{H}_{\varepsilon}) \cap \mathbb{R}\langle\varepsilon\rangle_b^n, \lim_{\varepsilon} \mathbf{x}_{\varepsilon} \neq \boldsymbol{\eta}\}$  is empty. Thus, the proof of Lemma 8.4.1 is finished.  $\square$

When  $\{\mathbf{y}_\varepsilon \in \text{crit}(\delta_{\mathbf{x}}, \mathcal{H}_\varepsilon) \cap \mathbb{R}\langle\varepsilon\rangle_b^n, \lim_\varepsilon \mathbf{y}_\varepsilon \neq \mathbf{x}\}$  is not empty, we prove the following criteria to identify whether the point  $\mathbf{x}$  is isolated in  $\mathcal{H}$ .

**Lemma 8.4.2.** *Let  $\mathbf{x}$  be a point of  $\mathcal{H} \cap \mathbb{R}^n$ ,  $\mathcal{C}_x$  be the connected component of  $\mathcal{H} \cap \mathbb{R}^n$  containing  $\mathbf{x}$  and  $\delta_{\mathbf{x}}$  is the function defined as above.*

*Assuming that*

$$\{\delta_{\mathbf{x}}(\lim_\varepsilon \mathbf{y}_\varepsilon) \mid \mathbf{y}_\varepsilon \in \text{crit}(\delta_{\mathbf{x}}, \mathcal{H}_\varepsilon) \cap \mathbb{R}\langle\varepsilon\rangle_b^n, \lim_\varepsilon \mathbf{y}_\varepsilon \neq \mathbf{x}\}$$

*is a non-empty finite set, we define a positive real number  $e_x$  as*

$$e_x = \min\{\delta_{\mathbf{x}}(\lim_\varepsilon \mathbf{y}_\varepsilon) \mid \mathbf{y}_\varepsilon \in \text{crit}(\delta_{\mathbf{x}}, \mathcal{H}_\varepsilon) \cap \mathbb{R}\langle\varepsilon\rangle_b^n, \lim_\varepsilon \mathbf{y}_\varepsilon \neq \mathbf{x}\}.$$

*Then, the following statements are equivalent:*

- i)  $\mathbf{x}$  is an isolated point of  $\mathcal{H} \cap \mathbb{R}^n$ .*
- ii) There exists  $e \in ]0, e_x[$  such that  $\mathcal{H} \cap S(\mathbf{x}, \sqrt{e}) = \emptyset$ .*
- iii) For every  $e \in ]0, e_x[$ ,  $\mathcal{H} \cap S(\mathbf{x}, \sqrt{e}) = \emptyset$ .*

*Moreover, if  $\mathbf{x}$  is not an isolated point of  $\mathcal{H} \cap \mathbb{R}^n$ , then  $\mathcal{C}_x$  intersects  $S(\mathbf{x}, \sqrt{e})$  for every  $e \in ]0, e_x[$ .*

*Proof.* By the definition of real isolated points, we immediately have that (i) implies (ii) and (iii) implies (i). It remains to demonstrate that (iii) is a consequence of (ii), which we separate into two statements: (ii) leads to (i) and then (i) leads to (iii).

We now show that (ii) implies (i) by contradiction. We assume that the point  $\mathbf{x}$  is not a real isolated point of  $\mathcal{H}$ . If  $\mathcal{C}_x$  is not bounded,  $\mathcal{C}_x$  intersects  $S(\mathbf{x}, \sqrt{e})$  for every  $e > 0$  and there is nothing to be proved. We now assume that  $\mathcal{C}_x$  is bounded.

By Lemma 8.2.2, there exists a point  $\mathbf{x}_\varepsilon$  such that  $\mathbf{x}_\varepsilon$  is bounded over  $\mathbb{R}$  and  $\lim_\varepsilon \mathbf{x}_\varepsilon = \mathbf{x}$ . Let  $\mathcal{C}_\varepsilon \subset \mathbb{R}\langle\varepsilon\rangle^n$  be a connected component of the real algebraic set  $\mathcal{H}_\varepsilon$  containing  $\mathbf{x}_\varepsilon$ . By Lemma 8.2.3,  $\mathcal{C}_\varepsilon$  is bounded over  $\mathbb{R}$ . Thus, by [9, Proposition 12.49], its limit is connected in  $\mathbb{R}^n$ . Consequently,  $\lim_\varepsilon \mathcal{C}_\varepsilon$  is a connected subset of  $\mathcal{C}_x$ .

Moreover, as  $\mathcal{C}_\varepsilon$  is closed and bounded,  $\delta_{\mathbf{x}}$  admits a maximum point  $\mathbf{y}_\varepsilon$  over  $\mathcal{C}_\varepsilon$  (see [19, Theorem 2.5.8]). Note that  $\mathbf{y}_\varepsilon \in \text{crit}(\delta_{\mathbf{x}}, \mathcal{H}_\varepsilon) \cap \mathbb{R}\langle\varepsilon\rangle^n$ . So, we have that  $\delta_{\mathbf{x}}(\lim_\varepsilon \mathbf{y}_\varepsilon) \geq e_x$ . Therefore, for any  $e \in ]0, e_x[$ , the closed ball  $\overline{B}(\mathbf{x}, \sqrt{e})$  does not contain  $\lim_\varepsilon \mathbf{y}_\varepsilon$ . Since  $\lim_\varepsilon \mathcal{C}_\varepsilon$  is connected in  $\mathbb{R}^n$  and contains  $\mathbf{x}$  and  $\lim_\varepsilon \mathbf{y}_\varepsilon$ , there exists a semi-algebraic continuous function  $\gamma : [0, 1] \rightarrow \lim_\varepsilon \mathcal{C}_\varepsilon$  such that  $\gamma(0) = \mathbf{x}$  and  $\gamma(1) = \lim_\varepsilon \mathbf{y}_\varepsilon$ . As  $\delta_{\mathbf{x}}(\mathbf{x}) = 0$  and  $\delta_{\mathbf{x}}(\lim_\varepsilon \mathbf{y}_\varepsilon) \geq e_x$ , by the intermediate value property [9, Proposition 3.5], there exists  $t_0 \in ]0, 1[$  such that  $\delta_{\mathbf{x}}(\gamma(t_0)) = e$  for any  $e \in ]0, e_x[$ . Therefore, the connected component  $\mathcal{C}_x$  intersects  $S(\mathbf{x}, \sqrt{e})$  at  $\gamma(t_0)$ .

So, (ii) does not hold either when  $\mathcal{C}_x$  is bounded. We conclude that (ii) leads to (i).

Finally, it remains to show that (i) implies (iii). Again, we prove this by contradiction. Assume that there exists  $e \in ]0, e_x[$  such that  $\mathcal{H} \cap S(\mathbf{x}, \sqrt{e}) \neq \emptyset$ . Equivalently, there exists a point  $\mathbf{z} \in \mathcal{H} \cap \mathbb{R}^n$  such that  $\delta(\mathbf{z}) = e \in ]0, e_x[$ .

By Lemma 8.2.2, there exists a point  $z_\varepsilon \in \mathcal{H}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^n$  such that  $\lim_\varepsilon z_\varepsilon = z$ . Let  $\mathcal{C}_{z,\varepsilon}$  be the connected component of  $\mathcal{H}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^n$  containing  $z_\varepsilon$ .

In the closed and connected semi-algebraic set  $\mathcal{C}_{z,\varepsilon}$ , there exists a point  $z'_\varepsilon$  at which the restriction of  $\delta_x$  to  $\mathcal{C}_{z,\varepsilon}$  reaches its minimum. So,  $z'_\varepsilon$  belongs to  $\text{crit}(\delta_x, \mathcal{H}_\varepsilon)$ . Thus, we have that

$$\delta_x(\lim_\varepsilon z'_\varepsilon) \leq \delta_x(\lim_\varepsilon z_\varepsilon) < e_x.$$

Using the definition of  $e_x$ , we deduce that  $\delta_x(\lim_\varepsilon z'_\varepsilon) = 0$ , which is equivalent to  $\lim_\varepsilon z'_\varepsilon = x$ .

So, both  $z_\varepsilon$  and  $z'_\varepsilon$  lie in the connected component  $\mathcal{C}_{z,\varepsilon}$  of  $\mathcal{H}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^n$  and  $\lim_\varepsilon z'_\varepsilon = x$ . If  $\mathcal{C}_{z,\varepsilon}$  is not bounded over  $\mathbb{R}$ , by Lemma 8.2.3, we have that  $\mathcal{C}_x$  is not bounded, which implies immediately that  $x$  is not an isolated point of  $\mathcal{H} \cap \mathbb{R}^n$ .

Otherwise, when  $\mathcal{C}_{z,\varepsilon}$  is bounded over  $\mathbb{R}$ , by [9, Proposition 12.49],  $\lim_\varepsilon \mathcal{C}_{z,\varepsilon}$  is a connected subset of  $\mathcal{H} \cap \mathbb{R}^n$  that contains  $x$  and  $z$ . In this case, we also conclude that  $x$  is not isolated in  $\mathcal{H} \cap \mathbb{R}^n$ . Therefore, (i) leads to (iii), which finishes our proof.  $\square$

For  $x \in \mathfrak{C} \cap \mathbb{R}^n$ , when  $\{\mathbf{y}_\varepsilon \in \text{crit}(\delta_x, \mathcal{H}_\varepsilon) \cap \mathbb{R}\langle\varepsilon\rangle_b^n, \lim_\varepsilon \mathbf{y}_\varepsilon \neq x\}$  is empty, we define

$$e_x = \min\{\delta_x(\lim_\varepsilon \mathbf{y}_\varepsilon) \mid \mathbf{y}_\varepsilon \in \text{crit}(\delta_x, \mathcal{H}_\varepsilon) \cap \mathbb{R}\langle\varepsilon\rangle_b^n, \lim_\varepsilon \mathbf{y}_\varepsilon \neq x\} = +\infty.$$

Let  $e_0 \in \mathbb{R}$  such that

$$0 < e_0 < e_x$$

for every  $x \in \mathfrak{C} \cap \mathbb{R}^n$ . We deduce from Lemmas 8.4.1 and 8.4.2 the following proposition, which is our main criteria for designing our second algorithm.

**Proposition 8.4.3.** *For any  $\mathbf{a} = (a_1, \dots, a_n) \in \mathbb{R}^n$  such that  $a_i > 0$  for  $1 \leq i \leq n$ , we define a value  $e_0 \in \mathbb{R}$  (depending on  $\mathbf{a}$ ) as above. Then, for any candidate  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_n) \in \mathfrak{C} \cap \mathbb{R}^n$ ,  $\boldsymbol{\eta}$  is an isolated point of  $\mathcal{H} \cap \mathbb{R}^n$  if and only if the following polynomial system*

$$f(x_1, \dots, x_n) = \sum_{i=1}^n a_i(x_i - \eta_i)^2 - e_0 = 0$$

*admits at least one solution in  $\mathbb{R}^n$ .*

*Proof.* We assume first that  $\{\mathbf{y}_\varepsilon \in \text{crit}(\delta_x, \mathcal{H}_\varepsilon) \cap \mathbb{R}\langle\varepsilon\rangle_b^n, \lim_\varepsilon \mathbf{y}_\varepsilon \neq x\}$  is not empty. Using Lemma 8.4.2, we have that  $\boldsymbol{\eta}$  is an isolated point of  $\mathcal{H}$  if and only if  $\mathcal{H}$  intersects the sphere  $S(\boldsymbol{\eta}, \sqrt{e_0})$ .

Otherwise, while  $\{\mathbf{y}_\varepsilon \in \text{crit}(\delta_x, \mathcal{H}_\varepsilon) \cap \mathbb{R}\langle\varepsilon\rangle_b^n, \lim_\varepsilon \mathbf{y}_\varepsilon \neq x\}$  is empty, by Lemma 8.4.1, then  $\mathcal{H}$  is either a single point  $\boldsymbol{\eta}$  or an unbounded connected component containing  $\boldsymbol{\eta}$ . The similar conclusion follows immediately, which ends our proof.  $\square$

## 8.4.2 Outline of the algorithm

In this subsection, we give the outline of an algorithm that identifies the real isolated points of  $\mathcal{H}$  based on the results of Subsection 8.4.1.

We start by choosing randomly from  $\mathbb{Q}_+^n$  a  $n$ -uple  $\mathbf{a} = (a_1, \dots, a_n)$ . Recall that Proposition 8.4.3 requires us to compute a value  $e_0 \in \mathbb{Q}$  such that

$$0 < e_0 < \min\{\delta_\eta(\lim_\varepsilon \mathbf{x}_\varepsilon) \mid \mathbf{x}_\varepsilon \in \text{crit}(\delta_\eta, \mathcal{H}_\varepsilon) \cap \mathbb{R}\langle\varepsilon\rangle^n, \lim_\varepsilon \mathbf{x}_\varepsilon \neq \boldsymbol{\eta}\},$$

for every  $\boldsymbol{\eta} \in \mathfrak{C} \cap \mathbb{R}^n$ , where  $\delta_\eta$  is defined as

$$(x_1, \dots, x_n) \mapsto a_1(x_1 - \eta_1)^2 + \dots + a_n(x_n - \eta_n)^2.$$

We call ComputeE0 a subroutine that takes as input the polynomial  $f \in \mathbb{Q}[x_1, \dots, x_n]$  and a  $n$ -uple  $\mathbf{a} \in \mathbb{Q}^n$  and returns such an  $e_0 \in \mathbb{Q}$ . The explicit description of ComputeE0 is given in Subsection 8.4.3.

With  $e_0$  obtained by ComputeE0, by Proposition 8.4.3, identifying whether a candidate  $\boldsymbol{\eta} \in \mathfrak{C} \cap \mathbb{R}^n$  is an isolated point of  $\mathcal{H} \cap \mathbb{R}^n$  can be done by deciding whether  $\mathcal{H} \cap S(\boldsymbol{\eta}, \sqrt{e_0}) = \emptyset$ .

Our algorithm will have the following outline in which Isolated is a subroutine that takes as input  $f, \mathcal{C}, \mathbf{a}$  and  $e_0$  and computes the isolating intervals  $\mathcal{B}$ .

---

### Algorithm 8.3: IsolatedPoints

---

**Input:** A polynomial  $f \in \mathbb{Q}[x_1, \dots, x_n]$   
**Output:** A zero-dimensional parametrization  $\mathcal{C}$  and a set  $\mathcal{B}$  of intervals of  $\mathbb{R}$

- 1  $A$  chosen randomly in  $\text{GL}(n, \mathbb{Q})$
- 2  $\mathcal{C} \leftarrow \text{Candidates}(f, A)$
- 3  $\mathbf{a}$  chosen randomly in  $\mathbb{Q}_+^n$
- 4  $e_0 \leftarrow \text{ComputeE0}(f, \mathcal{C}, \mathbf{a})$
- 5  $\mathcal{B} \leftarrow \text{Isolated}(f, \mathcal{C}, \mathbf{a}, e_0)$
- 6 **return**  $(\mathcal{C}, \mathcal{B})$

---

In what follows, we describe briefly two variants of Isolated; one allows us to obtain an algorithm whose arithmetic complexity lies in  $D^{O(n)}$  and the other is designed to obtain better practical performance. Both of these variants, in general, solve the polynomial system below over  $\mathbb{R}^n$  for each candidate  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_n) \in \mathfrak{C} \cap \mathbb{R}^n$ :

$$f(x_1, \dots, x_n) = a_1(x_1 - \eta_1)^2 + \dots + a_n(x_n - \eta_n)^2 - e_0 = 0. \quad (8.2)$$

However, since the candidates are encoded by the zero-dimensional parametrization  $\mathcal{C}$ , we cannot treat directly the system (8.2) and need some workarounds.

Assume that the zero-dimensional parametrization  $\mathcal{C}$  is given under polynomial form, i.e.,

$$Z(\mathcal{C}) = \{(v_1(t), \dots, v_n(t)) \mid w(t) = 0\}.$$

The first variant of `Isolated` considers the system

$$f(x_1, \dots, x_n) = a_1(x_1 - v_1(t))^2 + \dots + a_n(x_n - v_n(t))^2 - e_0 = w(t) = 0. \quad (8.3)$$

We now need to identify for which real roots of  $w(t)$ , the above system has at least one real solution. To do so, we basically compute a polynomial  $Q(t, z) \in \mathbb{Q}[t, z]$  such that  $Q(t_0, z)$  admits real solutions for the real root  $t_0$  of  $w(t)$  if and only if the system (8.3) has at least one solution with  $t = t_0$ . The problem is therefore reduced to a bivariate setting, which is easily solved by classical real root counting algorithms. The details for designing this subroutine is explained in Subsection 8.4.4.

For the second variant of `Isolated`, taking advantage of the knowledge of  $e_0$ , we replace the candidate  $\boldsymbol{\eta}$  in the system (8.2) by an “approximation”  $\tilde{\boldsymbol{\eta}}$  of  $\boldsymbol{\eta}$  whose coordinates lie in  $\mathbb{Q}$  and establish a similar result as Proposition 8.4.3 for these approximations (see Lemma 8.4.6).

Informally, we claim that a candidate  $\boldsymbol{\eta}$  is a real isolated point of  $\mathcal{H}$  if and only if the system

$$f(x_1, \dots, x_n) = a_1(x_1 - \tilde{\eta}_1)^2 + \dots + a_n(x_n - \tilde{\eta}_n)^2 - \frac{e_0}{4} = 0 \quad (8.4)$$

has no real solution. Therefore, once those approximations are identified, one can solve the system above over the reals using real root finding algorithms (see, e.g., [9, Chap. 16]).

The design of this variant is explained in Subsection 8.4.5.

### 8.4.3 Computing a value for $e_0$

In this subsection, we present an algorithm that computes a value  $e_0$  introduced in Proposition 8.4.3. This value allows us to remove the use of infinitesimals in the identification of isolated points later.

We recall that, for each  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_n) \in \mathbb{C}^n$ , the function  $\delta_{\boldsymbol{\eta}}$  is defined as

$$\begin{aligned} \delta_{\boldsymbol{\eta}} : \quad \mathbb{C}^n &\rightarrow \mathbb{C}, \\ \boldsymbol{x} = (x_1, \dots, x_n) &\mapsto \sum_{i=1}^n a_i(x_i - \eta_i)^2. \end{aligned}$$

To apply Lemma 8.4.2, we need to compute a value  $e_0 \in \mathbb{Q}$  such that

$$0 < e_0 < \min\{\delta_{\boldsymbol{\eta}}(\lim_{\varepsilon} \boldsymbol{x}_{\varepsilon}) \mid \boldsymbol{x}_{\varepsilon} \in \text{crit}(\delta_{\boldsymbol{\eta}}, \mathcal{H}_{\varepsilon}) \cap \mathbb{R}\langle \varepsilon \rangle_b^n, \lim_{\varepsilon} \boldsymbol{x}_{\varepsilon} \neq \boldsymbol{\eta}\}$$

for every  $\boldsymbol{\eta} \in \mathfrak{C} \cap \mathbb{R}^n$ . In the lemma below, we show that, for a generic choice of  $\boldsymbol{a}$ , every critical locus  $\text{crit}(\delta_{\boldsymbol{\eta}}, \mathcal{H}_{\varepsilon})$  is finite.

**Lemma 8.4.4.** *Let  $\boldsymbol{\eta} \in \mathfrak{C}$  be a candidate. Then there exists a non-empty Zariski open subset  $\mathcal{A}$  of  $\mathbb{C}^n$  such that, for  $\boldsymbol{a} = (a_1, \dots, a_n) \in \mathcal{A} \cap \mathbb{Q}_+^n$ , the critical locus of  $\delta_{\boldsymbol{\eta}}$  restricted to  $\mathcal{H}_{\varepsilon}$  is finite.*

*Proof.* Since  $f = \varepsilon$  defines a smooth algebraic subset  $V(f - \varepsilon)$  of  $\mathbb{C}\langle \varepsilon \rangle^n$ , the critical locus of the restriction of  $\delta_{\boldsymbol{\eta}}$  to  $V(f - \varepsilon)$

$$f - \varepsilon = y \cdot \frac{\partial f}{\partial x_1} - 2a_1(x_1 - \eta_1) = \dots = y \cdot \frac{\partial f}{\partial x_n} - 2a_n(x_n - \eta_n) = 0. \quad (8.5)$$

Now we consider  $\mathbf{a} = (a_1, \dots, a_n)$  and  $y$  as indeterminates. Let  $\varphi : \mathbb{C}^n \times \mathbb{C}^n \times \mathbb{C}$  be the polynomial mapping defined as

$$(\mathbf{x}, \mathbf{a}, y) \mapsto \left( f - \varepsilon, y \cdot \frac{\partial f}{\partial x_1} - 2a_1(x_1 - \eta_1), \dots, y \cdot \frac{\partial f}{\partial x_n} - 2a_n(x_n - \eta_n) \right).$$

Let  $\mathcal{X}$  be the non-empty Zariski open subset of  $\mathbb{C}$  defined as

$$(x_1 - \eta_1) \cdots (x_n - \eta_n) \neq 0.$$

The Jacobian matrix of  $\varphi$  with respect to  $(\mathbf{x}, \mathbf{a}, y)$

$$\begin{bmatrix} \frac{\partial f}{\partial x_1} & \cdots & \frac{\partial f}{\partial x_n} & 0 & 0 & \cdots & 0 \\ * & \cdots & * & \frac{\partial f}{\partial x_1} & -2(x_1 - \eta_1) & \cdots & 0 \\ \vdots & \ddots & * & \vdots & \vdots & \ddots & \vdots \\ * & \cdots & * & \frac{\partial f}{\partial x_n} & 0 & \cdots & -2(x_n - \eta_n) \end{bmatrix}$$

has full rank when  $\mathbf{x} \in \mathcal{X}$  and

$$f - \varepsilon = y \cdot \frac{\partial f}{\partial x_1} - 2a_1(x_1 - \eta_1) = \cdots = y \cdot \frac{\partial f}{\partial x_n} - 2a_n(x_n - \eta_n) = 0.$$

By Thom's weak transversality theorem (Theorem 2.5.12), there exists a non-empty Zariski open subset  $\mathcal{A}_\emptyset$  of  $\mathbb{C}^n$  such that, for  $\mathbf{a} \in \mathcal{A}_\emptyset$ ,  $\mathbf{0}$  is a regular value of the restriction of  $\varphi_{\mathbf{a}}$  to  $\mathcal{X}$ . Thus, for  $\mathbf{a} \in \mathcal{A}_\emptyset$ , by Jacobian criterion, the restriction of the solutions of

$$f - \varepsilon = y \cdot \frac{\partial f}{\partial x_1} - 2a_1(x_1 - \eta_1) = \cdots = y \cdot \frac{\partial f}{\partial x_n} - 2a_n(x_n - \eta_n) = 0$$

to  $\mathcal{X}$  is a finite set, i.e.,  $\text{crit}(\delta_{\boldsymbol{\eta}}, V(f - \varepsilon)) \cap \mathcal{X}$  is finite.

Now we study the restriction of  $\text{crit}(\delta_{\boldsymbol{\eta}}, V(f - \varepsilon))$  to  $\mathbb{C}^n \setminus \mathcal{X}$ . We choose  $\mathbf{a} \in \mathbb{Q}_+^n$ . Let  $I$  be a non-empty proper subset of  $\{1, \dots, n\}$  and  $\mathcal{X}_I$  be the subset of  $\mathbb{C}^n$  defined by

$$x_i = \eta_i \quad \text{for } i \in I \quad \text{and} \quad x_i \neq \eta_i \quad \text{for } i \notin \{1, \dots, n\} \setminus I.$$

Let  $\mathbf{x} \in \text{crit}(\delta_{\boldsymbol{\eta}}, V(f - \varepsilon)) \cap \mathcal{X}_I$ . As  $f(\boldsymbol{\eta}) = 0$ ,  $\boldsymbol{\eta} \notin V(f - \varepsilon)$  and  $\mathbf{x} \neq \boldsymbol{\eta}$ . Hence,  $y \neq 0$  in the system (8.5). Hence,  $y \cdot \frac{\partial f}{\partial x_i} - 2a_i(x_i - \eta_i) = 0$  implies that  $\frac{\partial f}{\partial x_i} = 0$ . Since  $V(f - \varepsilon)$  is smooth,  $\frac{\partial f}{\partial x_i}$  for  $i \in \{1, \dots, n\} \setminus I$  cannot vanish simultaneously at  $\mathbf{x}$ . This means

$$\{\mathbf{x} \mid \mathbf{x} \in \text{crit}(\delta_{\boldsymbol{\eta}}, V(f - \varepsilon)), x_i = \eta_i \text{ for } i \in I\}$$

coincides with the critical locus  $\text{crit}(\delta_{\boldsymbol{\eta}, I}, V(f_I - \varepsilon))$  where

$$\delta_{\boldsymbol{\eta}, I} : (x_j)_{j \in \{1, \dots, n\} \setminus I} \mapsto \sum_{j \in \{1, \dots, n\} \setminus I} a_j (x_j - \eta_j)^2$$

and  $f_I$  is the polynomial obtained from  $f$  by substituting  $x_i = \eta_i$  for  $i \in I$ . Therefore, we can use the same arguments as above to prove the following.

There exists a non-empty Zariski open subset  $\mathcal{A}_I$  of  $\mathbb{C}^n$  such that for  $\mathbf{a} \in \mathcal{A}_I \cap \mathbb{Q}_+^n$ , the restriction of  $\text{crit}(\delta_\eta, V(f - \varepsilon))$  to  $\mathcal{X}_I$  is finite.

Let  $\mathcal{A}_+ = \bigcap_{I \subseteq \{1, \dots, n\}} \mathcal{A}_I$  which is a non-empty Zariski open subset of  $\mathbb{C}^n$ . Given any  $\mathbf{a} \in \mathcal{A}_+ \cap \mathbb{Q}_+^n$ ,  $\text{crit}(\delta_\eta, V(f - \varepsilon))$  is a finite set. Similarly for  $V(f + \varepsilon)$ , we obtain a non-empty Zariski open subset  $\mathcal{A}_-$ . Taking the intersection  $\mathcal{A} = \mathcal{A}_+ \cap \mathcal{A}_-$  ends the proof.  $\square$

Since the set of candidates  $\mathfrak{C}$  is encoded by a zero-dimensional parametrization  $\mathcal{C}$ , we do the whole computation at once through the function  $\delta$  defined as

$$\begin{aligned} \delta : \quad \mathbb{C}^n \times \mathbb{C} &\rightarrow \mathbb{C}, \\ (x_1, \dots, x_n, t) &\mapsto \sum_{i=1}^n a_i (x_i - v_i(t))^2. \end{aligned}$$

The following lemma is immediate.

**Lemma 8.4.5.** *Let  $\mathcal{H}_{\varepsilon, t} \subset \mathbb{C}\langle \varepsilon \rangle^{n+1}$  be the algebraic set defined by  $f^2 - \varepsilon^2 = w(t) = 0$ . Then, the set of critical values  $\delta(\text{crit}(\delta, \mathcal{H}_{\varepsilon, t}))$  is the union of  $\delta_\eta(\text{crit}(\delta_\eta, \mathcal{H}_\varepsilon))$  for  $\eta \in \mathfrak{C}$ .*

*Proof.* The set  $\text{crit}(\delta, \mathcal{H}_{\varepsilon, t})$  are defined by the points of  $\mathcal{H}_{\varepsilon, t}$  at which the matrix

$$\begin{bmatrix} \frac{\partial f^2}{\partial x_1} & \cdots & \frac{\partial f^2}{\partial x_n} & 0 \\ \frac{\partial \delta}{\partial \delta} & \cdots & \frac{\partial \delta}{\partial x_n} & \frac{\partial \delta}{\partial t} \\ 0 & \cdots & 0 & w'(t) \end{bmatrix}$$

has rank at most 2.

As  $w(t)$  is square-free, for every  $t_0$  such that  $w(t_0) = 0$ ,  $w'(t_0)$  is not zero. Therefore, the condition above restricted to  $\mathcal{H}_{\varepsilon, t}$  is equivalent to

$$\text{rank} \begin{bmatrix} \frac{\partial f^2}{\partial x_1} & \cdots & \frac{\partial f^2}{\partial x_n} \\ \frac{\partial \delta}{\partial x_1} & \cdots & \frac{\partial \delta}{\partial x_n} \end{bmatrix} \leq 1.$$

For every complex root  $t_0$  of  $w(t)$ , let  $\boldsymbol{\eta}_0 = (v_1(t_0), \dots, v_n(t_0))$ . By fixing  $t = t_0$ , the rank deficiency above is reduced to

$$\text{rank} \begin{bmatrix} \frac{\partial f^2}{\partial x_1} & \cdots & \frac{\partial f^2}{\partial x_n} \\ \frac{\partial \delta_{\boldsymbol{\eta}_0}}{\partial x_1} & \cdots & \frac{\partial \delta_{\boldsymbol{\eta}_0}}{\partial x_n} \end{bmatrix} \leq 1,$$

which defines the set  $\text{crit}(\delta_{\boldsymbol{\eta}_0}, \mathcal{H}_\varepsilon)$ .

Thus,  $\text{crit}(\delta, \mathcal{H}_{\varepsilon, t}) = \bigcup_{w(t_0)=0} \{(\mathbf{x}, t_0) \mid \boldsymbol{\eta}_0 = (v_1(t_0), \dots, v_n(t_0)), \mathbf{x} \in \text{crit}(\delta_{\boldsymbol{\eta}_0}, \mathcal{H}_\varepsilon)\}$ . This concludes the proof.  $\square$

Now we aim to compute the limit of the critical points and their corresponding values of the restriction of  $\delta$  to the algebraic set  $\mathcal{H}_{\varepsilon,t}$  defined by  $f^2 - \varepsilon^2 = w(t) = 0$ . Note that Lemmas 8.4.4 and 8.4.5 imply that, for a generic  $\mathbf{a} \in \mathbb{Q}_+^n$ , the set of critical points  $\text{crit}(\delta, \mathcal{H}_{\varepsilon,t})$  is finite.

By [169, Theorems 1, 2], for generic values of  $\mathbf{a} \in \mathbb{Q}^n$ , we have that

$$\lim_{\varepsilon} \text{crit}(\delta, \mathcal{H}_{\varepsilon,t}) \subset \langle f \rangle + \left\langle w(t), y \cdot \frac{\partial f}{\partial x_i} - \frac{\partial \delta}{\partial x_i} \text{ for every } 1 \leq i \leq n \right\rangle \cap \mathbb{Q}[\mathbf{x}, t]$$

and the ideal on the right-hand side is zero-dimensional.

From the above inclusion, one can follow a similar computation as in Algorithm 8.1 using the geometric resolution algorithm [87]. However, as the degree of  $w(t)$  is bounded by  $2D(D-1)^{n-1}$  (see Subsection 8.3.4), such a computation would lead to an arithmetic complexity  $D^{O(n^2)}$ .

A workaround to obtain a better complexity is to use a variant of geometric resolution over the quotient ring  $\mathbb{A} = \mathbb{Q}[t]/\langle w(t) \rangle$  as explained in [174, Appendix J]. Note that  $w(t)$  is not necessarily irreducible, the extension  $\mathbb{A}$  is only a product of fields and doing the computation over the ring  $\mathbb{A}$  is not trivial. We will see in Subsection 8.4.6 that this approach allows us to obtain an algorithm with arithmetic complexity lying in  $D^{O(n)}$ .

Our subroutine ComputeE0 is designed as follows.

- a) First, we call a subroutine ParametricCurve that takes as input  $f, \mathcal{C}, \mathbf{a} \in \mathbb{Q}_+^n$  and  $i \in \{1, \dots, n\}$  and computes a one-dimensional parametrization  $\mathcal{J}_i$  over  $\mathbb{Q}[t]/\langle w(t) \rangle$  of the system

$$\left( \frac{\partial \delta}{\partial x_j} \cdot \frac{\partial f}{\partial x_i} - \frac{\partial \delta}{\partial x_i} \cdot \frac{\partial f}{\partial x_j} = 0 \right)_{j \in \{1, \dots, n\} \setminus \{i\}} \quad \text{and} \quad \frac{\partial f}{\partial x_i} \neq 0.$$

An explicit description of this subroutine can be found in [174, Appendix J.5].

- b) Next, we compute a zero-dimensional parametrization  $\mathcal{E}_i$  of the intersection of  $\mathcal{H} = V(f)$  with the sets of solutions defined by the parametrizations  $\mathcal{J}_i$  above.

This is done by calling a subroutine IntersectCurve on the input  $\mathcal{J}_i$  and  $f$ , which is described also in [174, Appendix J.5].

- c) We then call a subroutine Union that computes a zero-dimensional parametrization  $\mathcal{E}$  that defines  $\cup_{i=1}^n Z(\mathcal{E}_i)$ .
- d) Finally, taking as input the zero-dimensional parametrization  $\mathcal{E}$ , we call a subroutine GetE0 that computes the required value  $e_0$ . This can be done by calling FGLM algorithm [67] to compute a polynomial  $P(e)$  whose solutions encode the values  $e = \sum_{i=1}^n a_i(x_i - v_i(t))^2$  for  $\mathbf{x} \in Z(\mathcal{E})$  and  $w(t) = 0$ . Next, we evaluate a lower bound of the minimal distance between the roots of  $P(e)$  using [9, Proposition 10.23].

---

**Algorithm 8.4: ComputeE0**

---

**Input:**  $f \in \mathbb{Q}[x_1, \dots, x_n]$ ,  $\mathcal{C} = (w(t), v_1(t), \dots, v_n(t))$  and  $\mathbf{a} \in \mathbb{Q}_+^n$   
**Output:**  $e_0 \in \mathbb{Q}$   
1  $\delta \leftarrow a_1(x_1 - v_1(t))^2 + \dots + a_n(x_n - v_n(t))^2$   
2 **for**  $1 \leq i \leq n$  **do**  
3      $\mathcal{J}_i \leftarrow \text{ParametricCurve}(f, \mathbf{a}, \mathcal{C}, i)$   
4      $\mathcal{E}_i \leftarrow \text{IntersectCurve}(\mathcal{J}_i, f)$   
5  $\mathcal{E} \leftarrow \text{Union}(\mathcal{E}_1, \dots, \mathcal{E}_n)$   
6  $e_0 \leftarrow \text{GetE0}(\mathcal{E})$   
7 **return**  $e_0$

---

### 8.4.4 The first variant of IsIsolated

In this subsection, we explain the details of the first variant of IsIsolated.

Using the value  $e_0$  output by Algorithm 8.4, Proposition 8.4.3 allows one to identify the isolated points of  $\mathcal{H}$  among the candidates by checking whether the polynomial system

$$f(x_1, \dots, x_n) = \sum_{i=1}^n a_i(x_i - \eta_i)^2 - e_0 = 0$$

admits real solutions for each candidate  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_n) \in \mathfrak{C} \cap \mathbb{R}^n$ . Again, one can consider the system

$$f(x_1, \dots, x_n) = \sum_{i=1}^n a_i(x_i - v_i(t))^2 - e_0 = w(t) = 0 \quad (8.6)$$

to handle all the candidates at once.

Let  $\mathcal{W}_t \subset \mathbb{C}^{n+1}$  be the algebraic set defined the equation (8.6). Our strategy is to compute a finite subset of  $\mathcal{W}_t \cap \mathbb{R}^{n+1}$  that intersects every connected component of  $\mathcal{W}_t \cap \mathbb{R}^{n+1}$ . Then, all the real  $t$ -coordinates of those sample points correspond to the isolated points of  $\mathcal{H} \cap \mathbb{R}^n$ .

We consider the polynomial

$$F = f(x_1, \dots, x_n)^2 + \left( \sum_{i=1}^n a_i(x_i - v_i(t))^2 - e_0 \right)^2.$$

Note that  $F + w(t)^2$  defines also the real algebraic set  $\mathcal{W}_t \cap \mathbb{R}^{n+1}$ . Therefore, the sample points above can be computed using the algorithm of [169] on the input  $F + w(t)^2 \in \mathbb{Q}[t, x_1, \dots, x_n]$ . Such an algorithm returns a zero-dimensional parametrization over  $\mathbb{Q}$  that defines a finite set intersects every connected component of  $\mathcal{W}_t$ . Since the total degree of  $F + w(t)^2$  can go up to  $O(D^n)$ , this computation faces the same complexity issue as in Subsection 8.4.3. Again, we can bypass this problem by solving over  $\mathbb{A}[x_1, \dots, x_n]$  where  $\mathbb{A} = \mathbb{Q}[t]/\langle w(t) \rangle$ .

Let  $B$  be a matrix randomly chosen from  $\text{GL}(n, \mathbb{Q})$  and  $F^B(\mathbf{x}) = F(B \cdot \mathbf{x})$ . We apply the geometric resolution algorithm over  $\mathbb{A}$  on the system of equations:

$$F^B = 0, \quad \frac{\partial F^B}{\partial x_j} = 0, \quad \frac{\partial F^B}{\partial x_1} \neq 0.$$

This algorithm returns a zero-dimensional parametrization  $(U(z), V_1(z), \dots, V_n(z))$  over the ring  $\mathbb{A}$ , which means that  $V_1, \dots, V_n$  and  $U$  are elements of  $\mathbb{A}[z]$ , such that, for any real solution  $t_0$  of  $w(t)$ , the finite set defined by

$$\{(V_1(t_0, z), \dots, V_n(t_0, z)) \mid z \in \mathbb{R}, U(t_0, z) = 0\}$$

intersects every connected component of

$$f(x_1, \dots, x_n) = \sum_{i=1}^n a_i (x_i - v_i(t_0))^2 - e_0 = 0.$$

Hence, the isolated points of  $\mathcal{H} \cap \mathbb{R}^n$  are indeed

$$\{(v_1(t), \dots, v_n(t)) \mid (t, z) \in \mathbb{R}^2 : w(t) = U(t, z) = 0\}.$$

Our problem boils down to solving the bivariate system  $w(t) = U(t, z) = 0$  over  $\mathbb{R}^2$ .

In Algorithm 8.5 below, we introduce two subroutines:

- **BivariatePolynomial** takes as input the polynomials  $F$  and  $w(t)$  and returns the eliminating polynomial  $U(t, z)$ . It uses the geometric resolution algorithm over  $\mathbb{A}$  described in [174, Appendix J].
- **BivariateSolve** takes as input  $w(t)$  and  $U(t, z)$  and returns the set  $\mathcal{B}$  of intervals that isolate the real roots of  $w(t)$  corresponding to  $\mathcal{I}(\mathcal{H})$ . Such a subroutine can be designed efficiently with resultants.

---

**Algorithm 8.5:** The first variant of `Isolated`

---

**Input:**  $f, \mathcal{C} = (w(t), v_1(t), \dots, v_n(t)), \mathbf{a}$  and  $e_0$

**Output:** A set  $\mathcal{B}$  of isolating intervals

- 1  $F \leftarrow f(x_1, \dots, x_n)^2 + (\sum_{i=1}^n a_i (x_i - v_i(t))^2 - e_0)^2$
  - 2  $U(t, z) \leftarrow \text{BivariatePolynomial}(F, w(t))$
  - 3  $\mathcal{B} \leftarrow \text{BivariateSolve}(w(t), U(t, z))$
  - 4 **return**  $\mathcal{B}$
-

### 8.4.5 Approximations of the candidates

This subsection is devoted to the design of the second variant of `Isolated`. Note that this variant does not require solving polynomial systems in the quotient ring  $\mathbb{Q}[t]/\langle w(t) \rangle$ . Further, we will see that it is based mostly on isolating candidates from the zero-dimensional parametrization  $\mathcal{C}$ . Moreover, handling the candidates through the zero-dimensional parametrization  $\mathcal{C}$  means that we include every complex point of  $\mathcal{C}$  into the computation, which could be an overkill. The subroutine presented in what follows (see Algorithm 8.6) considers only the real points of  $\mathcal{C}$ .

The main idea is to replace the candidate  $\boldsymbol{\eta}$  in the criteria provided by Proposition 8.4.3 by a rational point  $\tilde{\boldsymbol{\eta}} \in \mathbb{Q}^n$ , which can be thought of an approximation of  $\boldsymbol{\eta}$ . To do so, we need to identify how close the points  $\boldsymbol{\eta}$  and  $\tilde{\boldsymbol{\eta}}$  need to be. The lemma below shows that requiring  $d_{\mathbf{a}}(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) < \sqrt{e_0}/2$  is enough.

**Lemma 8.4.6.** *Let  $\boldsymbol{\eta}$  be a candidate and  $\tilde{\boldsymbol{\eta}}$  be a point in  $\mathbb{R}^n$  satisfying  $d_{\mathbf{a}}(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) < \sqrt{e_0}/2$ . Then,  $\boldsymbol{\eta}$  is an isolated point of  $\mathcal{H} \cap \mathbb{R}^n$  if and only if  $\mathcal{H}$  does not intersect the sphere  $S(\tilde{\boldsymbol{\eta}}, \sqrt{e_0}/2)$ .*

*Proof.* If the set  $\{\mathbf{x}_\varepsilon \in \text{crit}(\delta_{\boldsymbol{\eta}} \cap \mathbb{R}\langle \varepsilon \rangle_b^n, \mathcal{H}_\varepsilon), \lim_\varepsilon \mathbf{x}_\varepsilon \neq \boldsymbol{\eta}\}$  is empty, then, by Lemma 8.4.1,  $\mathcal{H}$  is either a single point  $\boldsymbol{\eta}$  or an unbounded connected set containing  $\boldsymbol{\eta}$ . In either case, the conclusion of Lemma 8.4.6 is immediate. Thus, in what follows,  $\{\mathbf{x}_\varepsilon \in \text{crit}(\delta_{\boldsymbol{\eta}}, \mathcal{H}_\varepsilon) \cap \mathbb{R}\langle \varepsilon \rangle_b^n, \lim_\varepsilon \mathbf{x}_\varepsilon \neq \boldsymbol{\eta}\}$  is assumed to be non-empty.

We prove now the necessary implication. Assume that  $\boldsymbol{\eta}$  is an isolated point of  $\mathcal{H} \cap \mathbb{R}^n$ . By Lemma 8.4.2, the intersection of  $\mathcal{H}$  and  $S(\boldsymbol{\eta}, \sqrt{e})$  is empty for every  $e \in ]0, e_0[$ . So,  $\boldsymbol{\eta}$  is the only point of  $\mathcal{H}$  lying in the open ball  $B(\boldsymbol{\eta}, \sqrt{e_0})$ .

Since  $d_{\mathbf{a}}(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) < \sqrt{e_0}/2$ , the candidate  $\boldsymbol{\eta}$  does not lie on the sphere  $S(\tilde{\boldsymbol{\eta}}, \sqrt{e_0}/2)$ . Moreover,  $S(\tilde{\boldsymbol{\eta}}, \sqrt{e_0}/2)$  is contained in the open ball  $B(\boldsymbol{\eta}, \sqrt{e_0})$ . Then, we have that  $S(\tilde{\boldsymbol{\eta}}, \sqrt{e_0}/2) \cap \mathcal{H} = \emptyset$ .

Now we turn to the sufficient implication. Assume by contradiction that  $\boldsymbol{\eta}$  is not isolated in  $\mathcal{H} \cap \mathbb{R}^n$ . By Lemma 8.4.2, the connected component  $\mathcal{C}_{\boldsymbol{\eta}}$  of  $\mathcal{H}$  containing  $\boldsymbol{\eta}$  intersects the sphere  $S(\boldsymbol{\eta}, \sqrt{e_0})$ . So, there exists a semi-algebraic continuous function  $\gamma : [0, 1] \rightarrow \mathcal{C}_{\boldsymbol{\eta}}$  such that  $\gamma(0) = \boldsymbol{\eta}$  and  $\gamma(1)$  lying on the sphere  $S(\boldsymbol{\eta}, \sqrt{e_0})$ .

We have that

$$d_{\mathbf{a}}(\gamma(1), \tilde{\boldsymbol{\eta}}) \geq d_{\mathbf{a}}(\gamma(1), \boldsymbol{\eta}) - d_{\mathbf{a}}(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) > \sqrt{e_0} - \sqrt{e_0}/2 = \sqrt{e_0}/2.$$

As  $d_{\mathbf{a}}(\gamma(0), \tilde{\boldsymbol{\eta}}) < \sqrt{e_0}/2$  and  $d_{\mathbf{a}}(\gamma(1), \tilde{\boldsymbol{\eta}}) > \sqrt{e_0}/2$ , by the intermediate value property [9, Proposition 3.5], there exists  $t_0 \in ]0, 1[$  such that  $d_{\mathbf{a}}(\gamma(t_0), \tilde{\boldsymbol{\eta}}) = \sqrt{e_0}/2$ . This implies that the intersection of  $\mathcal{H}$  and  $S(\tilde{\boldsymbol{\eta}}, \sqrt{e_0}/2)$  is not empty, which concludes our proof.  $\square$

Let  $t_{\boldsymbol{\eta}}$  be the real root of  $w(t)$  that corresponds to  $\boldsymbol{\eta}$ , i.e.,  $\boldsymbol{\eta} = (v_1(t_{\boldsymbol{\eta}}), \dots, v_n(t_{\boldsymbol{\eta}}))$ . To apply Lemma 8.4.6, we need to choose  $t_{\tilde{\boldsymbol{\eta}}} \in \mathbb{Q}$  such that the rational point  $\tilde{\boldsymbol{\eta}} = (v_1(t_{\tilde{\boldsymbol{\eta}}}), \dots, v_n(t_{\tilde{\boldsymbol{\eta}}}))$  satisfies that  $d_{\mathbf{a}}(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) < \sqrt{e_0}/2$ . This leads us to identify  $\rho > 0$  such that  $|t_{\boldsymbol{\eta}} - t_{\tilde{\boldsymbol{\eta}}}| < \rho$  implies

$$a_1(v_1(t_{\boldsymbol{\eta}}) - v_1(t_{\tilde{\boldsymbol{\eta}}}))^2 + \dots + a_n(v_n(t_{\boldsymbol{\eta}}) - v_n(t_{\tilde{\boldsymbol{\eta}}}))^2 < \frac{e_0}{4}.$$

Lemma 8.4.7 below allows us to compute explicitly an appropriate value for  $\rho$ .

**Lemma 8.4.7.** Let  $\{t_1, \dots, t_\ell\}$  be the distinct real roots of  $w(t) = 0$  and  $\{\eta_1, \dots, \eta_\ell\}$  be the corresponding candidates.

We consider a set of intervals  $(I_j)_{1 \leq j \leq \ell}$  such that

- The intervals  $I_j$  are pairwise disjoint.
- The interval  $I_j$  contains only  $t_j$  as a real root of  $w(t)$ .

For each  $1 \leq i \leq n$ , let  $K_i = \max_{j=1}^{\ell} \max_{t \in I_j} |v'_i(t)|$ . Then, for any  $1 \leq j \leq \ell$  and  $t_\theta$  such that  $t_\theta \in I_j$  and  $|t_\theta - t_j| < \frac{1}{K_i} \cdot \sqrt{\frac{e_0}{4na_i}}$ , we have the following inequality:

$$|v_i(t_\theta) - v_i(t_j)| < \sqrt{\frac{e_0}{4na_i}}.$$

Let  $\rho \leq \min_{i=1}^n \frac{1}{K_i} \cdot \sqrt{\frac{e_0}{4na_i}}$ . For any real root  $t_\eta$  of  $w(t)$  and  $t_\theta \in I_j$  such that  $|t_\theta - t_\eta| < \rho$ , we have that

$$d_a(\theta, \eta) < \frac{\sqrt{e_0}}{2}.$$

*Proof.* For  $1 \leq j \leq \ell$  and any  $t_\theta \in \mathbb{Q}$ , we have that

$$v_i(t_\theta) - v_i(t_j) = v'_i(\tilde{t}_j)(t_\theta - t_j),$$

where  $\tilde{t}_j \in \mathbb{R}$  lies between  $t_\theta$  and  $t_j$ .

For  $t \in I_j = ]r_j, s_j[$ , by the definition of  $K_i$ , we have  $|v'_i(t)| \leq K_i$ . Then, for  $t_\theta \in I_j$  such that  $|t_\theta - t_j| < \frac{1}{K_i} \cdot \sqrt{\frac{e_0}{4na_i}}$ , we have

$$|v_i(t_\theta) - v_i(t_j)| = |v'_i(\tilde{t}_j) \cdot (t_\theta - t_j)| \leq K_i \cdot |t_\theta - t_j| < \sqrt{\frac{e_0}{4na_i}}.$$

Now we take  $\rho \leq \min_{i=1}^n \frac{1}{K_i} \cdot \sqrt{\frac{e_0}{4na_i}}$ . If  $t_\theta \in I_j$  and  $|t_\theta - t_j| < \rho$ , then we have

$$d_a(\theta, \eta_j) = \sqrt{\sum_{i=1}^n a_i (v_i(t_\theta) - v_i(t_j))^2} < \sqrt{\sum_{i=1}^n \frac{e_0}{4n}} = \frac{\sqrt{e_0}}{2}.$$

□

Lemmas 8.4.6 and 8.4.7 provides us the ingredients to design Algorithm 8.6. It requires us to introduce two subroutines *Isolate* and *MaxOverInterval* below.

- We need two versions of *Isolate*. The first one takes as input a polynomial  $p \in \mathbb{Q}[t]$  and returns a set of disjoint intervals of rational extremities isolating the real roots of  $p$ .

Besides the polynomial  $p \in \mathbb{Q}[t]$ , the second version of *Isolate* requires a positive  $\rho \in \mathbb{Q}$  as input and returns the intervals of length at most  $\rho$  that isolate the real roots of  $p$ .

The explicit descriptions of both of these real root isolating algorithms are given in [168].

- `MaxOverInterval` takes as input a polynomial  $p \in \mathbb{Q}[t]$  and an interval  $[r, s]$  where  $r, s \in \mathbb{Q}$  and returns an upper bound of  $\max_{t \in [r, s]} |p(t)|$ . Such a subroutine can be implemented using the following naive bound:

$$\max_{t \in [r, s]} |p(t)| \leq \sum_{i=0}^{\deg(p)} |c_i|$$

where  $p((s - r)t + r) = c_0 \cdot t^{\deg(p)} + \dots + c_{\deg(p)}$ .

Algorithm 8.6 proceeds through these following steps:

- We call `Isolate` on the input  $w(t)$  to obtain a set of intervals  $I_j$  that isolate the real roots of  $w(t)$  and compute  $K_i = \max_{j=1}^{\ell} \max_{t \in I_j} |v'_i(t)|$  using the subroutine `MaxOverInterval` on the input  $v'_i(t)$  and each interval  $I_j$ .
- We then compute  $\rho \in \mathbb{Q}$  such that  $0 < \rho \leq \min_{i=1}^n \frac{1}{K_i} \cdot \sqrt{\frac{e_0}{4na_i}}$  and use `Isolate` on the polynomial  $w(t)$  and the precision  $\rho$  to obtain a set of intervals  $\tilde{I}_j$  such that each  $\tilde{I}_j$  contains exactly one real root of  $w(t)$  and  $|\tilde{I}_j| < \rho$ .
- For  $1 \leq j \leq \ell$ , we choose a point  $\tilde{t}_j$  in  $I_j \cap \tilde{I}_j \cap \mathbb{Q}$  and evaluate  $\tilde{\eta}_j = (v_1(\tilde{t}_j), \dots, v_n(\tilde{t}_j))$ . The set  $\tilde{\mathcal{C}}$  of the approximations is taken as  $\{(\tilde{\eta}_j, I_j \mid 1 \leq j \leq \ell)\}$ .
- Finally, we decide whether the system

$$f(x_1, \dots, x_n) = \sum_{i=1}^n a_i(x_i - \tilde{\eta}_i)^2 - \frac{e_0}{4} = 0$$

has a real solution for each approximation  $\tilde{\eta}$  and return those which do not.

We summarize Section 8.4 in Algorithm 8.6 below, which is our second variant of `Isolated`.

---

**Algorithm 8.6:** Algorithm lIsolated-Approx

---

**Input:**  $f, \mathcal{C} = (w(t), v_1(t), \dots, v_n(t))$ ,  $\mathbf{a} \in \mathbb{Q}_+^n$  and  $e_0 \in \mathbb{Q}$   
**Output:** A set of isolating interval  $\mathcal{B}$

- 1  $\{I_1, \dots, I_\ell\} \leftarrow \text{Isolate}(w(t))$
- 2 **for**  $i \in \{1, \dots, n\}$  **do**
- 3      $K_i \leftarrow \max_{j=1}^{\ell} \text{MaxOverInterval}(v_i'(t), I_j)$
- 4  $\{\tilde{I}_1, \dots, \tilde{I}_\ell\} \leftarrow \text{Isolate}\left(w(t), \rho = \min_{i=1}^n \frac{1}{K_i} \cdot \sqrt{\frac{e_0}{4na_i}}\right)$
- 5 **for**  $j \in \{1, \dots, \ell\}$  **do**
- 6      $\tilde{t}_j \in I_j \cap \tilde{I}_j$
- 7      $\tilde{\eta}_j \leftarrow (v_1(\tilde{t}_j), \dots, v_n(\tilde{t}_j))$
- 8  $\tilde{\mathcal{C}} \leftarrow \{(\tilde{\eta}_j, I_j) \mid 1 \leq j \leq \ell\}$ ,  $\mathcal{B} \leftarrow \emptyset$
- 9 **for**  $(\tilde{\eta}, I_\eta) \in \tilde{\mathcal{C}}$  **do**
- 10     **if**  $\text{HasRealSolutions}(\tilde{\eta}, f, \mathbf{a}, e_0) = \text{false}$  **then**
- 11          $\mathcal{B} \leftarrow \mathcal{B} \cup I_\eta$
- 12 **return**  $\mathcal{B}$

---

**Remark for more efficient implementation.** From a computational point of view, checking the intersection of  $\mathcal{H} \cap \mathbb{R}^n$  with a sphere defined by a quadric would increase the bit-size coefficients appearing originally in  $f$ . Actually, we can take any hypercube such that it contains the candidate  $\eta$  in its interior and is contained in the ball  $B(\eta, e_0)$  and check whether the boundary of this hypercube intersects  $\mathcal{H} \cap \mathbb{R}^n$ . This leads us to check the emptiness of semi-algebraic sets defined by  $f = 0$  and some linear polynomial inequalities; the polynomials involved in such a computation have smaller degrees and bit-sizes than the ones for computing with the sphere.

### 8.4.6 Complexity analysis

The main objective of this subsection is to establish the complexity results for two variants of Algorithm 8.3, which use respectively Algorithm 8.5 and Algorithm 8.6 for the subroutine lIsolated. We start with the complexity estimate for Algorithm 8.5.

**Theorem 8.1.2.** *Let  $f \in \mathbb{Q}[x_1, \dots, x_n]$ . Then, the variant of Algorithm 8.3 which uses Algorithm 8.5 computes the real isolated points of the algebraic hypersurface defined by  $f$  within  $O^\sim(64^n D^{8n})$  arithmetic operations in  $\mathbb{Q}$  and one call of real root isolation on a univariate polynomial of degree bounded by  $2^{n+2} D^{2n}$ .*

*Proof.* Recall that, in Subsection 8.3.4, it is proved that computing the parametrization  $\mathcal{C}$  encoding the candidates can be done within  $O^\sim(D^{3n})$  arithmetic operations in  $\mathbb{Q}$  and the degrees of the polynomials  $w(t), v_1(t), \dots, v_n(t)$  are bounded by  $2D(D-1)^{n-1}$ . It remains to estimate the arithmetic complexity of the subroutines ComputeE0 and lIsolated.

Let  $\kappa$  be the degree of  $w(t)$ . Algorithm 8.4 (ComputeE0) relies on computing the limit of  $\text{crit}(\delta, \mathcal{H}_{\varepsilon,t}) \cap \mathbb{C}\langle\varepsilon\rangle_b^n$ , where  $\mathcal{H}_{\varepsilon,t}$  is the algebraic set defined by

$$(f - \varepsilon) \cdot (f + \varepsilon) = 0, \quad w(t) = 0 \quad (8.7)$$

and  $\delta$  is the distance function

$$(x_1, \dots, x_n) \mapsto \sum_{i=1}^n a_i (x_i - v_i(t))^2$$

We use the algorithm of [169] on the function  $\delta$  for the resolution of polynomial systems in the quotient ring  $\mathbb{Q}[t]/\langle w(t) \rangle$ . Using Bézout's bound on the system

$$f^2 - \varepsilon^2 = y \cdot \frac{\partial f}{\partial x_1} - \frac{\partial \delta}{\partial x_1} = \dots = y \cdot \frac{\partial f}{\partial x_n} - \frac{\partial \delta}{\partial x_n} = 0$$

defining  $\text{crit}(\delta, \mathcal{H}_{\varepsilon,t})$  over  $\mathbb{A}$ , the degree of  $\text{crit}(\delta, \mathcal{H}_{\varepsilon,t})$  in  $\mathbb{C}\langle\varepsilon\rangle^n$  is bounded by  $2D^{n+1}\kappa \leq 4D^{n+2}(D-1)^{n-1} \approx 4D^{2n+1}$ .

By [174, Appendix J.5], the arithmetic operations over  $\mathbb{A}$  can be done using  $O^{\sim}(\kappa)$  operations in  $\mathbb{Q}$ . Thus, applying [169, Theorem 5], we obtain the complexity bound  $O^{\sim}(\kappa \cdot D^{3n+2}) \approx O^{\sim}(D^{4n+2})$  for obtaining the zero-dimensional parametrization  $\mathcal{E}$  in Algorithm 8.4.

The call to GetE0 computes from the zero-dimensional parametrization  $\mathcal{E}$  a univariate polynomial  $P(e) \in \mathbb{Q}[e]$  whose solutions are the critical values of  $\delta$  restricted to  $\mathcal{V}_{\varepsilon}$ . Since the degree of  $P(e)$  is bounded by  $4D^{2n+1}$ , this can be done using FGLM algorithm [67] within  $O^{\sim}(D^{6n+3})$  arithmetic operations over  $\mathbb{Q}$ . Next, it computes the minimal distance between the real roots of  $P(e)$  using [9, Proposition 10.23]. The complexity of this computation is linear in the degree of  $P(E)$ . Thus, it does not change the asymptotic complexity of Algorithm 8.4.

Therefore, Algorithm 8.4 can be done within  $O^{\sim}(D^{6n+3})$  arithmetic operations in  $\mathbb{Q}$ .

Algorithm 8.5 is basically computing sample points of the hypersurface

$$F = f(x_1, \dots, x_n)^2 + \left( \sum_{i=1}^n a_i (x_i - v_i(t))^2 - e_0 \right)^2$$

over the quotient ring  $\mathbb{Q}[t]/\langle w(t) \rangle$ . Again, we follow the algorithm of [169] on the input  $F$  with the extended version of geometric resolution to the quotient ring  $\mathbb{A}$ . By [169, Theorem 6] with the overcost  $O^{\sim}(\kappa)$  of arithmetic operations over  $\mathbb{A}$ , we obtain the complexity bound

$$O^{\sim}(\kappa \cdot (2D)^{3n+2}) \approx O^{\sim}(8^n D^{4n+2})$$

for MinimalPolynomial.

The output polynomial  $U(t, z)$  has degree at most  $(2D)^n$  in  $z$  and  $\kappa$  in  $t$  so its total degree is bounded by  $(2D)^n + \kappa$ . Therefore, solving the bivariate system

$$w(t) = U(t, z) = 0$$

can be done within

$$O^\sim\left(\left((2D)^n + \kappa\right)^4 \kappa^2 \left((2D)^n + \kappa\right)^2\right) \approx O^\sim(64^n D^{8n})$$

arithmetic operations in  $\mathbb{Q}$  using geometric resolution. In the end, one needs to isolate the real roots of the eliminating polynomial output by the geometric resolution. That polynomial has degree bounded by  $\kappa((2D)^n + \kappa) \leq 2^{n+2} D^{2n}$ .

Adding up all the steps, we obtain the arithmetic complexity of Algorithm 8.3, which lies in  $O^\sim(64^n D^{8n})$  with a call to real root isolation on a polynomial of degree bounded by  $2^{n+2} D^{2n}$ .  $\square$

Note that for implementing our algorithm, we would mainly rely on Algorithm 8.6 (Isolated-Approx). Hence, we dedicate the rest of this subsection to discuss its complexity. The complexity result of our algorithm using Algorithm 8.6 is stated as follows.

**Theorem 8.1.3.** *Let  $f \in \mathbb{Q}[x_1, \dots, x_n]$ . Then, the variant of Algorithm 8.3 which uses Algorithm 8.6 requires  $O^\sim(D^{6n+3})$  arithmetic operations in  $\mathbb{Q}$  and two real root isolating calls on a univariate polynomial of degree bounded by  $2D(D-1)^{n-1}$ .*

*Proof.* Recall that Algorithm 8.6 computes an approximation  $\tilde{\eta} = (\tilde{\eta}_1, \dots, \tilde{\eta}_n)$  for each candidate  $\eta \in \mathfrak{C} \cap \mathbb{R}^n$  and decides whether the system

$$f(x_1, \dots, x_n) = \sum_{i=1}^n a_i (x_i - \tilde{\eta}_i)^2 - \frac{e_0}{4} = 0$$

has a real solution. The arithmetic complexity for solving each of those decision problems lies in  $O^\sim(8^n D^{3n+2})$  using [169]. Since the cardinality of  $\mathfrak{C}$  is bounded by  $2D(D-1)^{n-1}$ , Algorithm 8.6 runs within

$$O^\sim(8^n D^{4n+2})$$

arithmetic operations in  $\mathbb{Q}$ .

Note that all the complexities above are dominated by the complexity

$$O^\sim(D^{6n+3})$$

of ComputeE0 (Algorithm 8.4).

It remains to estimate the complexity of computing the approximations, whose main steps consist of calling MaxOverInterval and isolating the real roots of the eliminating polynomial  $w(t)$  in the zero-dimensional parametrization encoding  $\mathfrak{C}$ .

The subroutine MaxOverInterval is called  $n$  times for the polynomials  $v'_i(t)$ ; this would require  $O^\sim(\deg(w)) \approx O^\sim(D^n)$  arithmetic operations over  $\mathbb{Q}$ .

Since each of  $\ell$  evaluations  $\tilde{\eta}_j \leftarrow (v_1(\tilde{t}_j), \dots, v_n(\tilde{t}_j))$  takes  $O(nD^n)$  arithmetic operations, the cost of getting the approximations is bounded by

$$O(nD^{2n}).$$

The real root isolation is called twice in Algorithm 8.6 on the polynomial  $w(t)$ .

Summing up the above discussion, we conclude that Algorithm 8.6 requires

$$O \sim (8^n D^{4n+2})$$

arithmetic operations in  $\mathbb{Q}$  and two calls of real root isolation on a univariate polynomial of degree bounded by  $2D(D-1)^{n-1}$ .  $\square$

Furthermore, the complexity of real root isolation algorithms depends on the degree of the input polynomial, which is  $w(t)$  in our case, and its bit-size of coefficients. For instance, using the algorithm of [176], we obtain a bit complexity

$$O \sim (\deg(w)^3 \tau^2)$$

where  $\tau$  is the largest bit-size of coefficients of  $w(t)$ . While the degree of  $w(t)$  is already bounded by  $2D(D-1)^{n-1}$ ,  $\tau$  is not estimated yet in this thesis. To identify a bound of  $\tau$ , one needs to estimate the bit complexity of Algorithm 8.1 (Candidates), especially the algorithm for computing at least one point per connected component of a real algebraic set given in [169]. This topic will be studied in future research.

## 8.5 Optimizations

Even though computing the constant  $e_0$  requires at most  $D^{O(n)}$  arithmetic operations in  $\mathbb{Q}$ , its performance depends heavily on an efficient implementation of the geometric resolution algorithm over  $\mathbb{Q}[t]/\langle w(t) \rangle$ , which remains challenging to obtain. Thus, we aim to avoid such a computation as much as possible. In what follows, we present two subroutines which are launched to test whether it is necessary for computing  $e_0$ . In most of the case, with these subroutines, our algorithm will return the set of isolated points without doing any further computation.

### 8.5.1 Simple identification of real isolated points

The optimization subroutine described in what follows computes efficiently a subset of the candidates whose elements are real isolated points of  $\mathcal{H}$ . To do this, we identify for each candidate  $\mathbf{x}$  a ball  $B \in \mathbb{R}^n$  such that when the intersection of the boundary of  $B$  and  $\mathcal{H} \cap \mathbb{R}^n$  is empty,  $\mathbf{x}$  is isolated in  $\mathcal{H} \cap \mathbb{R}^n$  (but not the inverse). This allows us to avoid the computation of  $e_0$ .

We start with defining the set

$$\mathfrak{C}_2 = \bigcup_{i=1}^n \lim_{\varepsilon} \text{crit}(\pi_i, \mathcal{H}_\varepsilon) \cap \mathfrak{C}\langle \varepsilon \rangle_b^n.$$

Note that the set of candidates  $\mathfrak{C}$  is a subset of  $\mathfrak{C}_2$ . We have the following lemmas.

**Lemma 8.5.1.** *For every bounded connected component  $\mathcal{C}$  of  $\mathcal{H} \cap \mathbb{R}^n$  that is not a singleton, there exist at least two points in  $\mathfrak{C}_2$  that belong to  $\mathcal{C}$ .*

*Proof.* Let  $\mathcal{C}_1, \dots, \mathcal{C}_k$  be the connected components of  $\mathcal{H}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^n$  such that  $\lim_\varepsilon \mathcal{C}_i \subset \mathcal{C}$ . By Lemma 8.2.3, since  $\mathcal{C}$  is bounded, the  $\mathcal{C}_i$ 's are bounded over  $\mathbb{R}$ . Then, by [9, Proposition 12.49],  $\lim_\varepsilon \mathcal{C}_i$  is connected. On the other hand,  $\bigcup_{i=1}^k \lim_\varepsilon \mathcal{C}_i = \mathcal{C}$ . As  $\mathcal{C}$  is not a singleton, there exists  $1 \leq i \leq k$  such that  $\lim_\varepsilon \mathcal{C}_i$  is not a singleton either.

Now, since  $\mathcal{C}_i$  is a connected component of  $\mathcal{H}_\varepsilon \cap \mathbb{R}\langle\varepsilon\rangle^n$  that is bounded over  $\mathbb{R}$  and not a singleton. Then, there exists a coordinate  $x_j$  such that the projection of  $\lim_\varepsilon \mathcal{C}_i$  on  $x_j$  is an interval which is not a point. Consequently, the projection of  $\mathcal{C}_i$  on the  $x_j$ -coordinate is a closed interval  $[\alpha, \beta] \subset \mathbb{R}\langle\varepsilon\rangle$ . Then, there exist two points  $\mathbf{x}_\alpha$  and  $\mathbf{x}_\beta$  such that  $\mathbf{x}_\alpha$  and  $\mathbf{x}_\beta$  are in  $\text{crit}(\pi_j, \mathcal{H}_\varepsilon) \cap \mathcal{C}_i$  and  $\pi(\mathbf{x}_\alpha) = \alpha$  and  $\pi(\mathbf{x}_\beta) = \beta$ . Since  $\pi_j(\lim_\varepsilon \mathbf{x}_\alpha) \neq \pi_j(\lim_\varepsilon \mathbf{x}_\beta)$ ,  $\lim_\varepsilon \mathbf{x}_\alpha \neq \lim_\varepsilon \mathbf{x}_\beta$ . Then  $\lim_\varepsilon \mathbf{x}_\alpha$  and  $\lim_\varepsilon \mathbf{x}_\beta$  are two distinct points of  $\mathcal{C} \cap \lim_\varepsilon \text{crit}(\pi_j, \mathcal{H}_\varepsilon)$ . As a consequence, there exists two distinct points in  $\mathfrak{C}_2$  that lie on  $\mathcal{C}$ .  $\square$

**Proposition 8.5.2.** *Let  $\mathfrak{C}_2$  be defined as above. Let  $\mathbf{x} \in \mathfrak{C} \cap \mathbb{R}^n$  and  $B \subset \mathbb{R}^n$  be a ball such that  $\mathfrak{C}_2 \cap B = \{\mathbf{x}\}$  and  $\mathbf{x}$  is contained in the interior of  $B$ . Then, if the intersection of the boundary of  $B$  and  $\mathcal{H} \cap \mathbb{R}^n$  is empty,  $\mathbf{x}$  is an isolated point of  $\mathcal{H} \cap \mathbb{R}^n$ .*

*Proof.* Let  $\mathcal{C}$  be the connected component of  $\mathcal{H} \cap \mathbb{R}^n$  containing  $\mathbf{x}$ . If  $\mathcal{C}$  is unbounded, then  $\mathbf{x}$  is not an isolated point and the intersection of the boundary  $B$  with  $\mathcal{H} \cap \mathbb{R}^n$  is not empty. We now assume that  $\mathcal{C}$  is bounded.

We assume by contradiction that  $\mathcal{C}$  is not a singleton. By Lemma 8.5.1, there exists a point  $\mathbf{y} \in \mathfrak{C}_2 \cap \mathbb{R}^n$  such that  $\mathbf{y} \neq \mathbf{x}$  and  $\mathbf{y} \in \mathcal{C}$ . Since  $\mathbf{x}$  and  $\mathbf{y}$  lie on different sides of  $B$ . Hence, by intermediate value theorem, the intersection of the boundary of  $B$  and  $\mathcal{H} \cap \mathbb{R}^n$  is not empty, which ends the proof.  $\square$

Recall that, in the subroutine Candidates, the zero-dimensional parametrization encoding  $\lim_\varepsilon \text{crit}(\pi_i, \mathcal{H}_\varepsilon) \cap \mathbb{R}\langle\varepsilon\rangle_b^n$  are already computed. Therefore, the union  $\mathfrak{C}_2$  can be obtained easily by taking the union of those zero-dimensional parametrizations. Next, we isolate the candidates in  $\mathfrak{C} \cap \mathbb{R}^n$  by balls such that each of them contains exactly one point of  $\mathfrak{C}_2 \cap \mathbb{R}^n$ .

Algorithm 8.7 contains the description of the subroutine SimpleIdentification. We call to a subroutine BoxIsolate that takes as input a zero-dimensional parametrization encoding a subset of  $\mathbb{C}^n$  and computes isolating boxes for its real zeros.

---

**Algorithm 8.7: SimpleIdentification**

---

**Input:** A zero-dimensional parametrization  $\mathcal{C}_2$

**Output:** A set  $\mathcal{B}_1$  of intervals of  $\mathbb{R}$

```
1  $\mathcal{B}_1 \leftarrow \emptyset$ 
2 Boxes  $\leftarrow$  BoxIsolate( $\mathcal{C}_2$ )
3 for box  $\in$  Boxes do
4   if box  $\cap \mathcal{H} = \emptyset$  then
5      $\mathcal{B}_1 \leftarrow \mathcal{B}_1 \cup \{t\text{-coordinate of box}\}$ 
6 return  $\mathcal{B}_1$ 
```

---

By Proposition 8.5.2, for each  $\mathbf{x} \in \mathfrak{C}$  such that the intersection of the ball isolating  $\mathbf{x}$  with  $\mathcal{H}$  is empty, we conclude that  $\mathbf{x}$  is an isolated point of  $\mathcal{H}$ . For the non-empty intersections, we cannot decide whether  $\mathbf{x}$  is isolated yet. The problem arises when the isolating boxes are not small enough so that they intersect not only the connected component of  $\mathcal{H} \cap \mathbb{R}^n$  containing  $\mathbf{x}$  but also some other connected component. When this happens, one could try a smaller size of isolating boxes.

### 8.5.2 Limits of critical curves

To compute a set of candidates, we consider the critical loci  $\text{crit}(\pi_i, \mathcal{H}_\varepsilon)$  for  $1 \leq i \leq n$ . Our second optimization considers the critical loci of the projections on the plane; especially, the limits of those critical loci are curves in  $\mathbb{R}^n$  whose real isolated points contain the isolated points of  $\mathcal{H} \cap \mathbb{R}^n$ . Thus, one can compute a superset of  $\mathcal{I}(\mathcal{H})$  through computing the real isolated points of limits of critical curves.

More precisely, for  $1 \leq i < j \leq n$ , we denote by  $\pi_{i,j}$  the projection

$$\pi_{i,j} : (x_1, \dots, x_n) \mapsto (x_i, x_j).$$

Recall that  $\mathcal{H}_\varepsilon$  is a smooth algebraic set defined by

$$(f - \varepsilon) \cdot (f + \varepsilon) = 0.$$

**Lemma 8.5.3.** *Let  $A \in \text{GL}(n, \mathbb{Q})$ . For every  $1 \leq i < j \leq n$ , the set of isolated points of  $\mathcal{H}^A \cap \mathbb{R}^n$  is contained in set of real isolated points of  $\lim_\varepsilon \text{crit}(\pi_{i,j}, \mathcal{H}_\varepsilon^A)$ .*

*Proof.* Let  $\mathbf{x}$  be an isolated point of  $\mathcal{H}^A \cap \mathbb{R}^n$ . By Proposition 8.2.4,  $\mathbf{x} \in \lim_\varepsilon \text{crit}(\pi_i, \mathcal{H}_\varepsilon^A)$ . Since  $\text{crit}(\pi_i, \mathcal{H}_\varepsilon^A) \subset \text{crit}(\pi_{i,j}, \mathcal{H}_\varepsilon^A)$ , we obtain  $\mathbf{x} \in \lim_\varepsilon \text{crit}(\pi_{i,j}, \mathcal{H}_\varepsilon^A)$ . Note that  $\lim_\varepsilon \text{crit}(\pi_{i,j}, \mathcal{H}_\varepsilon^A)$  is a subset of  $\mathcal{H}^A$ . Thus, if  $\mathbf{x}$  is isolated in  $\mathcal{H}^A \cap \mathbb{R}^n$ , it is also isolated in  $\lim_\varepsilon \text{crit}(\pi_{i,j}, \mathcal{H}_\varepsilon^A) \cap \mathbb{R}^n$ .  $\square$

**Remark 8.5.4.** *Note that a real isolated point of  $\lim_\varepsilon \text{crit}(\pi_{i,j}, \mathcal{H}_\varepsilon)$  is not necessarily isolated in  $\mathcal{H} \cap \mathbb{R}^n$ . Take for example the degenerate torus, given by the equation*

$$(x_1^2 + x_2^2 + x_3)^2 - 4(x_1^2 + x_2^2) = 0.$$

The real trace of  $\lim_{\varepsilon} \text{crit}(\pi_{1,2}, \mathcal{H}_{\varepsilon})$  is the union of the point  $(0, 0)$  and the circle given by

$$x_1^2 + x_2^2 - 4 = x_3 = 0.$$

Hence, Lemma 8.5.3 allows us to obtain a superset of  $\mathcal{I}(\mathcal{H})$  only.

By [171, Theorem 2], for a generic change of variables  $A$ , the critical locus  $\text{crit}(\pi_{i,j}, \mathcal{H}_{\varepsilon}^A)$  is an equidimensional algebraic set of dimension one defined by

$$(f - \varepsilon) \cdot (f + \varepsilon) = 0, \quad \frac{\partial f}{\partial x_k} = 0 \quad \text{for } 1 \leq k \leq n \text{ and } k \neq i, j.$$

The computation of  $\lim_{\varepsilon} \text{crit}(\pi_{i,j}, \mathcal{H}_{\varepsilon}^A)$  can be done using a similar subroutine of Subsection 8.2.2. For each  $1 \leq i, j \leq n$ , we denote by  $J_{i,j}$  the ideal

$$\left\langle \frac{\partial f}{\partial x_k} = 0 \text{ for } 1 \leq k \leq n, k \neq i, j \right\rangle.$$

**Lemma 8.5.5.** *Let  $\pi_{i,j}$  be defined as above. There exists a non-empty Zariski open subset  $\mathcal{A}$  of  $\text{GL}(n, \mathbb{C})$  such that, for any  $A \in \mathcal{A} \cap \text{GL}(n, \mathbb{Q})$ , the algebraic set  $\mathcal{C}$  defined by*

$$V \left( \langle f^A \rangle + J_k : \left( \frac{\partial f^A}{\partial x_i} \right)^{\infty} \cap J_k : \left( \frac{\partial f^A}{\partial x_j} \right)^{\infty} \right)$$

*is equi-dimensional of dimension 1 and contains  $\lim_{\varepsilon} \text{crit}(\pi_{i,j}, \mathcal{H}_{\varepsilon}^A)$ .*

*As a consequence, any isolated point of  $\mathcal{H}^A \cap \mathbb{R}^n$  is also isolated in  $\mathcal{C} \cap \mathbb{R}^n$ .*

*Proof.* The proof of the first statement follows a similar outline of the proof of [169, Theorem 1 and Theorem 2]. From the inclusion

$$\mathcal{I}(\mathcal{H}^A) \subset \lim_{\varepsilon} \text{crit}(\pi_{i,j}, \mathcal{H}_{\varepsilon}^A) \subset \mathcal{C} \cap \mathbb{R}^n \subset \mathcal{H}^A \cap \mathbb{R}^n,$$

we deduce that every real isolated point of  $\mathcal{H}^A \cap \mathbb{R}^n$  is also an isolated point of  $\mathcal{C} \cap \mathbb{R}^n$ .  $\square$

We define the subroutine CurveLimitCheck that takes as input  $f \in \mathbb{Q}[x_1, \dots, x_n]$  and  $A \in \text{GL}(n, \mathbb{Q})$  and returns a set of isolating boxes  $\mathcal{B}_2$ . It calls to two subroutines:

- CurveLimit that takes as input  $f$ ,  $A$  and a pair of index  $(i, j)$  and returns the eliminating polynomial of a rational parametrization encoding  $\lim_{\varepsilon} \text{crit}(\pi_{i,j}, \mathcal{H}_{\varepsilon}^A)$ . The design of this subroutine follows Lemma 8.5.5.
- Bivariatsolated that takes as input a bivariate polynomial  $U_{i,j}$  and computes the boxes isolating the real isolated points of  $V(U_{i,j})$ . This can be done by computing a cylindrical algebraic decomposition adapted to  $U_{i,j} = 0$ .

---

**Algorithm 8.8: CurveLimitCheck**

---

**Input:**  $f \in \mathbb{Q}[x_1, \dots, x_n]$ ,  $A \in \text{GL}(n, \mathbb{Q})$   
**Output:** A set  $\mathcal{B}_2$  of intervals of  $\mathbb{R}$

- 1 **for**  $1 \leq i < j \leq n$  **do**
- 2      $U_{i,j} \leftarrow \text{CurveLimit}(f, A, (i, j))$
- 3      $\text{boxes}_{i,j} \leftarrow \text{BivariateIsolated}(U_{i,j})$
- 4  $\mathcal{B}_2 \leftarrow \bigcap_{1 \leq i, j \leq n} \text{boxes}_{i,j}$
- 5 **return**  $\mathcal{B}_2$

---

**Summary.** To conclude this section, we show below the pseudo-code of our implementation. The subroutine `Candidates` is modified so that it returns, besides  $\mathcal{C}$  encoding the candidates, a zero-dimensional parametrization  $\mathcal{C}_2$  encoding the finite set  $\mathfrak{C}_2$ .

---

**Algorithm 8.9: Implementation of IsolatedPoints**

---

**Input:** A polynomial  $f \in \mathbb{Q}[x_1, \dots, x_n]$   
**Output:** A zero-dimensional parametrization  $\mathcal{C}$  and a set  $\mathcal{B}$  of intervals of  $\mathbb{R}$

- 1  $A$  chosen randomly in  $\text{GL}(n, \mathbb{Q})$
- 2  $\mathcal{C}, \mathcal{C}_2 \leftarrow \text{Candidates}(f, A)$
- 3  $\mathcal{B}_1 \leftarrow \text{SimpleIdentification}(\mathcal{C}_2)$
- 4 **if**  $|\mathcal{B}_1| = |\mathcal{C} \cap \mathbb{R}^n|$  **then**
- 5     **return**  $\mathcal{B}_1$
- 6  $\mathcal{B}_2 \leftarrow \text{CurveLimitCheck}(f, A)$
- 7 **if**  $|\mathcal{B}_1| = |\mathcal{B}_2|$  **then**
- 8     **return**  $\mathcal{B}_1$
- 9  $\mathbf{a}$  chosen randomly in  $\mathbb{Q}_+^n$
- 10  $e_0 \leftarrow \text{ComputeE0}(f, \mathcal{C}, \mathbf{a})$
- 11  $\mathcal{B} \leftarrow \text{Isolated-Approx}(f, \mathcal{C}, \mathbf{a}, e_0)$
- 12 **return**  $(\mathcal{C}, \mathcal{B})$

---

## 8.6 Experimental results

In this section, we report on practical performances of our algorithms. Computations were done on an Intel(R) Xeon(R) CPU E7-4820 2GHz and 1.5 TB of RAM. We take sums of squares of  $n$  random dense quadrics in  $n$  variables (with a non-empty intersection over  $\mathbb{R}$ ); we obtain *dense quartics* defining a finite set of points. Timings are given in seconds (s.), minutes (m.), hours (h.) and days (d.).

Table 8.1 shows the timings of Algorithm 8.2. Timings for Algorithm 8.1 (`Candidates`) are given in the column `CAND` below. Timings for the computation of the roadmaps are given in the column `RMP` and timings for the analysis of connectivity queries are given in the column `QRI`.

We use FGB library for computing Gröbner bases in order to perform algebraic elimination in our algorithms. We also used our C implementation for bivariate polynomial system solving (based on resultant computations) which we need to analyze connectivity queries in roadmaps.

Roadmaps are obtained as the union of critical loci of some maps in slices of the input variety [174]. We report on the highest degree of these critical loci in the column SRMP. The column SQRI reports on the maximum degree of the bivariate zero-dimensional system we need to study to analyze connectivity queries on the roadmap.

None of the examples we considered could be tackled using the implementation of CAD algorithm in Maple within 10 days.

We also implemented [9, Alg. 12.16] using the FLINT C library with evaluation/interpolation techniques instead to tackle coefficients involving infinitesimals. This algorithm only computes sample points per connected components. *That implementation was not able to compute sample points of the input quartics for any of our examples.* We then report in the column [BPR] on the number of complex solutions of the zero-dimensional system which is expected to be solved by [9, BPR]. This is to be compared with the columns SRMP and SQRI.

$n$	CAND	RMP	QRI	total	SRMP	SQRI	[BPR]
4	1 s.	15 s.	33 s.	50 s.	36	359	7290
5	20 s.	1h.	7h.	8 h.	108	4644	65 610
6	30 m.	2 d.	18 d.	20 d.	308	47952	590 490

Figure 8.1: Timings of Algorithm 8.2.

Table 8.2 below reports on the timings of Algorithm 8.9. Our implementation uses FGB library to perform algebraic elimination in our algorithms. The subroutine HasRealSolutions in Algorithm 8.6 is done by RAGLIB. Solving of zero-dimensional systems in the whole algorithm is done by MSOLVE and real root isolation is done by the command ROOTFINDING[ISOLATE] in Maple.

The column CAND2 shows the timings for computing the zero-dimensional parametrization  $\mathcal{C}_2$  and isolates its zeros (see Subsection 8.5.1). The column |real sols. /  $\deg(w)$  shows the number of real candidates among the total number of candidates. This motivates the use of approximations, which runs only on candidates in  $\mathbb{R}^n$ , instead of computing over  $\mathbb{Q}[t] / \langle w(t) \rangle$  which takes into account all candidates.

The column TEST1 reports on the timings of the first optimization (Algorithm 8.7). Exploiting the fact that isolating boxes are given by linear inequalities, we tweak RAGLIB for solving the associated decision problems. As explained in the end of Subsection 8.5.1, by isolating  $\mathcal{C}_2$  with a small enough boxes in the subroutine SimpleIdentification, one can also obtain a certified output without computing  $e_0$ . In our examples, it is the case and we do not need to carry out further computations. Timings of other steps are given as an indication for further researches.

Timings for ComputeE0 are given in the column E0. The columns APPROX and RAGLIB respectively give the timings for computing the approximations and solving the decision problem by RAGLIB. Note that the implementation used for two columns APPROX and RAGLIB checks the

emptiness of intersections of  $\mathcal{H} \cap \mathbb{R}^n$  with hypercubes (as explained in the end of Subsection 8.4.5). This computation is similar to the one of TEST1 with the main difference coming from the fact that the isolating boxes computed in APPROX requires more precision. This results in linear polynomials, that define hypercubes, of larger bit-sizes, which makes the column RAGLIB slower than TEST1.

At the moment, we do not dispose of a geometric resolution algorithm for  $\mathbb{Q}[t]/\langle w(t) \rangle$ . The implementation of ComputeE0 relies on available tools such as FGB, MSOLVE that work over the rational numbers only. The complexity of this subroutine is actually bounded by  $D^{O(n^2)}$  and the timings show that it is not practical. The value  $e_0$  in these examples is obtained since we know in advance that the real algebraic set is finite.

$n$	CAND	CAND2	real sols.   / $\deg(w)$	TEST1	<b>total</b>	E0	APPRX	RAGLIB
2	.1 s	.1 s	1/4	.1 s	.3 s	3 s	.1 s	.1 s
3	.2 s	.3 s	4/8	6 s	7 s	1 m	.1 s	10 s
4	1 s	4 s	2/16	1 m	1 m	20 h	.1 s	2 m
5	20 s	90 s	2/32	10 m	12 m	> 10 d	.2 s	15 m
6	30 m	2.5 h	2/64	4 h	7 h	> 10 d	20 s	6 h

Figure 8.2: Timings of Algorithm 8.9.

# Chapter 9

## Topics for future research

In the short and medium term, my research plan is to improve and generalize the results obtained in this thesis and to apply them to practical problems. Some research directions for each problem are presented in what follows.

### 9.1 Resolution of parametric polynomial systems

Concerning parametric polynomial systems, we can develop the results we obtained for real root classification and one block quantifier elimination in many aspects.

**Generalization to semi-algebraic sets.** The algorithms we presented in Chapters 5 and 6 are restricted to the case of polynomial systems of equations. One next step is to study how to deal with polynomial inequalities. Note that an inequality, strict or not, can be reformulated as an equation by using an extra variable

$$\begin{aligned}g \geq 0 &\rightarrow g - z^2 = 0, \\g > 0 &\rightarrow z^2 \cdot g - 1 = 0.\end{aligned}$$

This approach allows us to apply straight-forwardly our algorithms in Chapters 5 and 6 to handle inequalities. However, this method increases the number of variables and also the degree of the system taken as input. In practice, this affects the complexity of the algorithms and usually leads to inefficient computation.

When the input system involves one inequality, one can still rely on the classical theory of Hermite quadratic forms. More precisely, let  $I$  be a zero-dimensional ideal of  $\mathbb{Q}[x_1, \dots, x_n]$  and  $g \in \mathbb{Q}[x_1, \dots, x_n]$ . The quadratic form defined by

$$H(I, g) : \begin{array}{ccc} \mathbb{Q}[x_1, \dots, x_n]/I \times \mathbb{Q}[x_1, \dots, x_n]/I & \rightarrow & \mathbb{Q}, \\ (p, q) & \mapsto & \text{trace}(\mathcal{L}_{p \cdot q \cdot g}) \end{array}$$

allows one to compute the Tarski's query of  $I$  and  $g$  (Definition 4.4.3):

$$\text{TarskiQuery}(I, g) = |\{\eta \in V(I) \cap \mathbb{R}^n \mid g(\eta) > 0\}| - |\{\eta \in V(I) \cap \mathbb{R}^n \mid g(\eta) < 0\}|.$$

To obtain only  $|\{\eta \in V(I) \cap \mathbb{R}^n \mid g(\eta) > 0\}|$  or  $|\{\eta \in V(I) \cap \mathbb{R}^n \mid g(\eta) < 0\}|$ , one can combine  $H(I, 1)$ ,  $H(I, g)$  and  $H(I, g^2)$  that provide respectively their sum and difference. Hence, counting the number of points of  $V(I) \cap \mathbb{R}^n$  at which  $g > 0$  or  $g < 0$  requires three Hermite matrices.

Using a similar construction over the ground field  $\mathbb{Q}(\mathbf{y})$ , we can naturally extend this definition to parametric polynomial systems which contain one inequality.

When the given system contains multiple inequalities  $g_1, \dots, g_m$ , one can consider  $3^m$  Hermite matrices  $H(I, g_1^{\sigma_1} \cdots g_m^{\sigma_m})$  where  $\sigma_i \in \{0, 1, 2\}$ . Those Hermite matrices allow us to identify all the realizable sign conditions of  $(g_1, \dots, g_m)$ . However, such an approach would lead to an exponential complexity in  $m$ , which is not satisfying since the state-of-the-art complexity of one block quantifier elimination is only polynomial in the number of polynomials defining the input formula.

In [9, Chap. 10], a workaround is presented to avoid the exponential number of Tarski's query required for zero-dimensional ideals without parameters. It uses the so-called *adapted family* which detects and removes sign conditions that are not realizable along the computation. We would like to investigate further this direction to apply it to our algorithms for solving parametric systems with polynomial inequalities.

**The structure of Hermite matrices.** As observed in Chapter 5, we know that a parametric Hermite matrix with respect to the lexicographic ordering is a Hankel matrix. Therefore, such a matrix of size  $\delta \times \delta$  can be defined by only  $2\delta - 1$  elements.

Using grevlex ordering, we obtain a symmetric matrix with entries of smaller degrees but drop the structure of Hankel matrices. Even though, the number of distinct entries of those matrices can be much smaller than  $\frac{\delta(\delta+1)}{2}$ . Actually, for generic systems of  $n$  polynomials in  $\mathbb{Q}[x_1, \dots, x_n]$  of degree 2, we have the following table:

$n$	Distinct entries	$2^{n-1}(2^n + 1)$	Ratio
2	9	10	90%
3	27	36	75%
4	78	136	57%
5	224	528	42%

This structure depends on the staircase of the Gröbner basis with respect to the grevlex ordering of the input system. Therefore, a careful study of the combinatorics of those staircases would help to accelerate the computation with Hermite matrices, for instance, computing their determinants or signatures. We can rely on the results of the structure of staircases of generic systems [155, 159] or determinantal systems [43, 16] to study the structure of Hermite matrices.

Another research direction is to look at the geometric meaning of Hermite matrices. Recall that the rank of a Hermite matrix is equal to the number of distinct complex solutions of a zero-dimensional system. For a parametric polynomial system, the rank deficiency of its associated parametric Hermite matrix naturally leads to a partition of the parameter space.

To be precise, let  $\mathbf{f} = (f_1, \dots, f_s) \subset \mathbb{Q}[\mathbf{x}, \mathbf{y}]$  and  $H$  be the parametric Hermite matrix of  $\mathbf{f}$  with respect to the grevlex ordering. We denote the size of  $H$  by  $\delta$ . One can decompose the

parameter space  $\mathbb{C}^t$  into

$$\mathbb{C}^t = \bigcup_{i=1}^{\delta} (D_i \setminus D_{i-1})$$

where

$$D_r = \{\eta \in \mathbb{C}^t \mid \text{rank } \mathcal{H}(\eta) \leq r\}.$$

The set  $D_r$  is actually an algebraic set defined by  $(r+1)$ -minors of  $H$ . Moreover, when  $\mathbf{f}$  is taken to be a randomly generated dense system, we observe that  $D_r$  has dimension  $r$  for  $\delta-t+1 \leq r \leq \delta$  and  $D_r = \emptyset$  for  $r \leq \delta-t$ . We should emphasize that this behavior is different from generic matrices as in the example below.

**Example 9.1.1.** *We consider the polynomial  $f = x^3 + y_1x^2 + y_2x + y_3$ . The parametric Hermite matrix of  $\langle f \rangle$  with respect to the basis  $\{1, x, x^2\}$  is the Hankel matrix*

$$H = \begin{bmatrix} 3 & -y_1 & y_1^2 - 2y_2 \\ -y_1 & y_1^2 - 2y_2 & -y_1^3 + 3y_1y_2 - 3y_3 \\ y_1^2 - 2y_2 & -y_1^3 + 3y_1y_2 - 3y_3 & y_1^4 - 4y_1^2y_2 + 4y_1y_3 + 2y_2^2 \end{bmatrix}.$$

The locus at which the rank of  $H$  is at most 1 is an algebraic set of codimension 2 defined by  $\langle y_1^2 - 3y_2, y_1y_2 - 9y_3 \rangle$ .

On the other hand, such rank deficiency of a Hankel matrix  $M = (m_{i,j})_{1 \leq i,j \leq 3}$  where the  $m_{i,j}$ 's are indeterminates leads to an algebraic set of codimension 3 (see, e.g., [104]).

Recall that the most costly step of our algorithms for real root classification and one block quantifier elimination is the computation of sample points on the non-vanishing locus of some minors of Hermite matrices. Such a computation boils down to the resolution of zero-dimensional systems encoding some critical locus (see Section 5.2).

Since the complexity of solving zero-dimensional systems is polynomial in the number of solutions, decomposing these systems into new systems with smaller numbers of solutions would accelerate our algorithm. Therefore, we expect to be able to improve the computation of sample points by exploiting the decomposition of  $\mathbb{C}^t$  as above.

**Complexity results for non-generic inputs.** Recall that the complexity results provided in Chapters 5 and 6 rely on the genericity of the input  $\mathbf{f}$ , namely the Noether position of the homogenization of  $\mathbf{f}$ . Particularly, under this genericity assumption, the entries of the Hermite matrices are elements in  $\mathbb{Q}[\mathbf{y}]$  and we obtain a bound on the degrees of these entries (see Section 5.6).

However, in general, the entries of parametric Hermite matrices are rational functions in  $\mathbf{y}$ . They contain denominators coming from the leading coefficients (with respect to the  $\mathbf{x}$  variables) in Gröbner bases used in their construction. We illustrate this by the following toy example.

**Example 9.1.2.** *We consider the system*

$$\mathbf{f} = \{2x_1^2 + 2x_1x_2 + 2x_1y + 4x_2 + 2y + 1, 3x_1^2 + 4x_2y + y^2 + x_1 + y\}.$$

The parametric Hermite matrix of  $\mathbf{f}$  associated to the basis  $\{1, x_2, x_1\}$

$$\begin{bmatrix} 3 & \frac{-43y^2+49y-30}{12y} & \frac{4y-7}{3} \\ \frac{-43y^2+49y-30}{12y} & \frac{1225y^4-3566y^3+4933y^2-3384y+900}{144y^2} & \frac{-92y^3+241y^2-248y+90}{18y} \\ \frac{4y-7}{3} & \frac{-92y^3+241y^2-248y+90}{18y} & \frac{34y^2-62y+37}{9} \end{bmatrix}.$$

Studying the behavior of these denominators would help to control their degrees and could also allow us to tweak the algorithm to remove them. This research direction is therefore essential to obtain a complexity result and better practical performance for our algorithm in non-generic cases.

Let  $\mathbf{f}$  be a polynomial sequence in  $\mathbb{Q}[\mathbf{x}, \mathbf{y}]$ ,  $G$  be the reduced Gröbner basis of  $\mathbf{f}$  with respect to the ordering  $\text{grevlex}(\mathbf{x} \succ \mathbf{y})$  and  $H$  be the parametric Hermite matrix of  $\mathbf{f}$  associated to  $G$ .

Recall that the non-specialization polynomial of  $H$  is defined as

$$\mathbf{w}_\infty = \prod_{g \in G} \text{lc}_x(g).$$

Since the denominators in  $H$  appear as products of certain  $\text{lc}_x(g)$ , a reasonable approach to bound the degrees of denominators in  $H$  should start with controlling the degrees of those leading coefficients.

For instance, let  $g \in G$ . As  $\mathbb{Q}(\mathbf{y})$  is considered as the ground field, we can actually replace  $g$  by  $g/\text{lc}_x(g)$  in  $G$ . By substituting formally each  $\text{lc}_x(g)^{-1}$  by a new parameter  $z_g$ , one obtains polynomials in  $\mathbb{Q}[\mathbf{x}, \mathbf{y}, \mathbf{z}]$  where  $\mathbf{z} = \{z_g \mid g \in G\}$ . Using similar techniques as the ones in Section 5.6, we may be able to control the degrees of  $(\mathbf{y}, \mathbf{z})$  during the normal form reductions (which operate mainly on the variables  $\mathbf{x}$ ) in the construction of Hermite matrices.

The approach above could allow us to obtain a complexity result parameterized by the degrees of  $\text{lc}_x(g)$  for non-generic systems. Hence, estimating the degree of those leading coefficients  $\text{lc}_x(g)$  is an important question to be investigated.

In this research direction, we can also investigate input systems equipped with special structures such as symmetry or sparsity. Designing algorithms which are tailored for structured systems is important in computer algebra.

## 9.2 Total real intersection by hyperplanes

**Computing totally real hyperplane sections with higher multiplicities.** At this moment, using our real root classification algorithm, we are able to answer the (non-)existence of simple totally real hyperplane sections for multiple real algebraic curves. That is to decide for a real curve  $X$  whether there exists a hyperplane  $H$  defined over  $\mathbb{R}$  such that the intersection  $H \cap X$  contains only real points of *multiplicity* 1. These results lead to certain bounds of the simple real divisor bound for those examples. The next step is therefore to extend our algorithm to compute totally real hyperplane sections in general, i.e., the intersections with higher multiplicities.

Constructing a Hermite matrix associated to the parametric system under study, one can require that the rank of the Hermite matrix is equal to its signature, which is to say that the numbers of distinct complex and real solutions of that system are equal. This condition is then modeled by semi-algebraic formulas whose atoms are minors of the Hermite matrix. Hence, deciding the existence of totally real hyperplane sections is reduced to checking the emptiness of semi-algebraic sets defined by those formulas.

For instance, given a parametric Hermite matrix  $\mathcal{H}$  constructed as explained in Section 7.3 and  $r \in \mathbb{N}$ , a real point of parameters at which the  $(r + 1)$ -minors of  $\mathcal{H}$  vanish and the first  $r$  leading principal minors are positive would give a totally real hyperplane section of  $r$  points. This problem of deciding the emptiness leads to computations which are out-of-reach of classical tools and current libraries.

For instance, computing a general totally real hyperplane section of the curve  $X'_2$  in Example 7.4.1 would require to decide the emptiness of a semi-algebraic set in  $\mathbb{R}^3$  defined by one equation of degree 18 and 4 inequalities of degree 8, 10, 16 and 18. For the curve  $X_1$  of Example 7.4.1, we need to carry out similar computations on a semi-algebraic set in  $\mathbb{R}^3$  defined by an equation of degree 18 and 4 inequalities of degree 10, 12, 16 and 22.

We aim to design new algorithms that exploit the determinantal structure of those semi-algebraic sets in order to carry out those computations.

**Determining the deformation value.** Recall that, in Example 7.4.2, we revisit the counterexamples for Huisman's conjecture given in [125]. In these counterexamples, they construct a family of curves  $X_\varepsilon \in \mathbb{P}^3$  parameterized by  $\varepsilon > 0$  using deformation such that for a small enough  $\varepsilon > 0$ ,  $X_\varepsilon$  admits no totally real hyperplane section. Clearly, obtaining an explicit curve without totally real hyperplane section from the family  $X_\varepsilon$  requires an appropriate value of  $\varepsilon > 0$  which we want to identify automatically.

More specifically, we raise the computational problem to determine a number  $\varepsilon_0 \in \mathbb{R}_+$  such that any  $0 < \varepsilon < \varepsilon_0$  will lead to a curve  $X_\varepsilon$  without totally real hyperplane section. Representing  $\varepsilon$  by a new variable  $E$ , the family of curves  $X_\varepsilon$  is modeled by a surface in  $\mathbb{R}^4$  with coordinates  $(x_1, x_2, x_3, E)$ . Taking also the parameterized hyperplane

$$y_1x_1 + y_2x_2 + y_3x_3 + 1 = 0$$

into the defining system of  $X_\varepsilon$ , we obtain a system of 7 indeterminates which defines an algebraic set that we name  $\mathfrak{E}$  in  $\mathbb{C}^7$ .

Let  $\pi_E$  be the projection  $(\mathbf{x}, \mathbf{y}, E) \mapsto E$ . We want to compute a value  $\varepsilon_0$  such that for any  $\varepsilon \in (0, \varepsilon_0)$ , the fiber  $\pi_E^{-1}(\varepsilon) \cap \mathfrak{E} \cap \mathbb{R}^7$  contains at least one real point. Such a computation can be done by classifying the real solutions of the above system.

Another approach for solving this computational problem is to go through the generalized critical values of the restriction of  $\pi_E$  to  $\mathfrak{E} \cap \mathbb{R}^7$  (see, e.g., [127] for this definition). The knowledge of those generalized critical values would allow us to identify an  $\varepsilon_0 > 0$  such that there exists a

homeomorphism

$$(\pi_E^{-1}(\varepsilon) \cap \mathfrak{E} \cap \mathbb{R}^7) \times (0, \varepsilon_0) \simeq \pi_E((0, \varepsilon_0)) \cap \mathfrak{E} \cap \mathbb{R}^7$$

for every  $\varepsilon \in (0, \varepsilon_0)$ . Consequently, this value  $\varepsilon_0$  obtained from the generalized critical values provides the deformation value for our problem.

We want to investigate further the computation of  $\varepsilon_0$  by both approaches presented above.

### 9.3 Computing real isolated points

**Extension to semi-algebraic sets.** The algorithms presented in Chapter 5 allow us to compute the isolated points of a given real algebraic set. Recall that those algorithms take as input a single equation. When the given real algebraic set is known as the real solutions of a system of multiple equations

$$f_1 = \dots = f_s = 0,$$

one needs to do the computation through the sum of squares  $f_1^2 + \dots + f_s^2$ . However, doing so may increase the dimension of the underlying complex algebraic set and therefore makes the computation more difficult.

Therefore, it will be nice if our algorithms can be adapted to deal directly with systems of polynomial equations. Such an algorithm would require a variant of [169] that computes at least one point per connected component for a singular algebraic set defined by multiple equations.

Furthermore, to compute isolated points of semi-algebraic sets in general, one would need an algorithm that computes points per connected component of a given semi-algebraic set. Even though those algorithms exist, they handle the singular semi-algebraic sets through the deformation technique using infinitesimals. This would lead to inefficient implementation which cannot be used for solving applications. On the other hand, we currently reformulate the inequality constraints as equations with new variables. This direct method could lead to non-regular systems for which the computation of Gröbner bases is not practical.

As we aim to tackle applications in material sciences, we are interested in searching for a better approach to deal with inequalities and designing a practical algorithm for computing points per connected component of a singular semi-algebraic set.

**Bit-complexity of our algorithms.** In the implementation (Algorithm 8.9) of our algorithm for computing the real isolated points presented in Section 8.5, we make use of approximations of the candidates to avoid partly the computation over  $\mathbb{Q}[t]/\langle w(t) \rangle$ . Moreover, we introduce the optimization subroutine SimpleIdentification (Algorithm 8.7); this subroutine is able to replace the remaining computation of  $\sigma$  in many cases. These subroutines lead to an algorithm with better practical performance. However, they both rely on isolating the real solutions of the univariate polynomial  $w(t)$  in the zero-dimensional parametrization encoding the candidates. Whereas, the complexity of real root isolating depends on the bit-size of the coefficients of the eliminating

polynomial  $w(t)$  which is not yet analyzed in this thesis. Therefore, we want to estimate the bit complexity of our algorithms in the future. For this direction of research, we can refer to [57] for the bit complexity of the algorithm that computes sample points in [171].

# Bibliography

- [1] H. Anai, S. Hara, and K. Yokoyama. Sum of roots with positive real parts. In *Proceedings of the 2005 International Symposium on Symbolic and Algebraic Computation*, ISSAC '05, page 21–28, New York, NY, USA, 2005. Association for Computing Machinery.
- [2] H. Anai and V. Weispfenning. Reach set computations using real quantifier elimination. In *Hybrid Systems: Computation and Control*, pages 63–76. Springer Berlin Heidelberg, 2001.
- [3] D. S. Arnon. A cluster-based Cylindrical Algebraic Decomposition algorithm. *Journal of Symbolic Computation*, 5(1/2):189–212, 1988.
- [4] P. Aubry, D. Lazard, and M. Moreno Maza. On the theories of triangular sets. *Journal of Symbolic Computation*, 28(1):105–124, 1999.
- [5] A. Bardet. *Diviseurs sur les courbes réelles*. PhD thesis, Université d'Angers, 2013.
- [6] M. Bardet. *Étude des systèmes algébriques surdéterminés. Applications aux codes correcteurs et à la cryptographie*. Theses, Université Pierre et Marie Curie - Paris VI, Dec. 2004.
- [7] M. Bardet, J.-C. Faugère, and B. Salvy. On the complexity of the F5 Gröbner basis algorithm. *Journal of Symbolic Computation*, 70:49–70, 2015.
- [8] S. Basu, R. Pollack, and M.-F. Roy. On the combinatorial and algebraic complexity of quantifier elimination. *J. ACM*, 43(6):1002–1045, Nov. 1996.
- [9] S. Basu, R. Pollack, and M.-F. Roy. *Algorithms in Real Algebraic Geometry*. Springer-Verlag, Berlin, Heidelberg, 2006.
- [10] S. Basu and M. Roy. Divide and conquer roadmap for algebraic sets. *Discrete & Computational Geometry*, 52(2):278–343, 2014.
- [11] S. Basu, M. Roy, M. Safey El Din, and É. Schost. A baby step-giant step roadmap algorithm for general algebraic sets. *Foundations of Computational Mathematics*, 14(6):1117–1172, 2014.
- [12] D. Bayer and D. Mumford. What can be computed in Algebraic Geometry? pages 1–56, 1992.
- [13] D. Bayer and M. Stillman. A theorem on refining division orders by the reverse lexicographic order. *Duke Math. J.*, 55(2):321–328, 06 1987.
- [14] E. Becker, T. Mora, M. G. Marinari, and C. Traverso. The shape of the Shape lemma. In *Proceedings of the International Symposium on Symbolic and Algebraic Computation*, ISSAC '94, page 129–133, New York, NY, USA, 1994. Association for Computing Machinery.

- [15] C. Beltrán and A. Leykin. Robust certified numerical homotopy tracking. *Foundations of Computational Mathematics*, 13:253–295, 2013.
- [16] J. Berthomieu, A. Bostan, A. Ferguson, and M. Safey El Din. Gröbner bases and critical values: The asymptotic combinatorics of determinantal systems. preprint, 4 2021.
- [17] J. Berthomieu, C. Eder, and M. Safey El Din. msolve: A Library for Solving Polynomial Systems. Preprint, Feb. 2021.
- [18] C. Bick, M. Goodfellow, C. R. Laing, and E. A. Martens. Understanding the dynamics of biological and neural oscillator networks through exact mean-field reductions: a review. *Journal of mathematical neuroscience*, 10:9–9, 2020.
- [19] J. Bochnak, M. Coste, and M.-F. Roy. *Real Algebraic Geometry*. Springer-Verlag, Berlin, Heidelberg, 1998.
- [20] B. Bonnard, M. Chyba, A. Jacquemard, and J. Marriott. Algebraic geometric classification of the singular flow in the contrast imaging problem in nuclear magnetic resonance. *Mathematical Control & Related Fields*, 3(4):397–432, 2013.
- [21] B. Bonnard, J.-C. Faugère, A. Jacquemard, M. Safey El Din, and T. Verron. Determinantal sets, singularities and application to optimal control in medical imagery. In *Proceedings of the ACM on International Symposium on Symbolic and Algebraic Computation*, pages 103–110, 2016.
- [22] C. Borcea and I. Streinu. Geometric auxetics. *Proc. R. Soc. Lond., A, Math. Phys. Eng. Sci.*, 471(2184):24, 2015. Id/No 20150033.
- [23] C. S. Borcea and I. Streinu. Periodic auxetics: structure and design. *Q. J. Mech. Appl. Math.*, 71(2):125–138, 2018.
- [24] C. W. Brown. Improved projection for Cylindrical Algebraic Decomposition. *Journal of Symbolic Computation*, 32(5):447 – 465, 2001.
- [25] C. W. Brown. Qepcad b: A program for computing with semi-algebraic sets using cads. *SIGSAM Bull.*, 37(4):97–108, Dec. 2003.
- [26] C. W. Brown and J. H. Davenport. The complexity of quantifier elimination and cylindrical algebraic decomposition. In *Proceedings of the 2007 International Symposium on Symbolic and Algebraic Computation, ISSAC '07*, page 54–60, New York, NY, USA, 2007. Association for Computing Machinery.
- [27] C. W. Brown, M. El Kahoui, D. Novotni, and A. Weber. Algorithmic methods for investigating equilibria in epidemic modeling. *Journal of Symbolic Computation*, 41(11):1157–1173, 2006. Special Issue on the Occasion of Volker Weispfenning’s 60th Birthday.

- [28] W. Bruns and U. Vetter. *Determinantal Rings*. Lecture Notes in Mathematics. Springer-Verlag, Berlin, Heidelberg, 1988.
- [29] B. Buchberger. *Ein Algorithmus zum Auffinden der Basiselemente des Restklassenringes nach einem nulldimensionalen Polynomideal*. PhD thesis, University of Innsbruck, 1965.
- [30] B. Buchberger. A theoretical basis for the reduction of polynomials to canonical forms. *SIGSAM Bull.*, 10(3):19–29, Aug. 1976.
- [31] B. Buchberger. A criterion for detecting unnecessary reductions in the construction of Groebner bases. In *Proceedings of the International Symposium on Symbolic and Algebraic Computation*, EUROSAM '79, page 3–21, Berlin, Heidelberg, 1979. Springer-Verlag.
- [32] J. Cabral, H. Luckhoo, M. Woolrich, M. Joensson, H. Mohseni, A. Baker, M. L. Kringelbach, and G. Deco. Exploring mechanisms of spontaneous functional connectivity in meg: How delayed network interactions lead to structured amplitude envelopes of band-pass filtered oscillations. *NeuroImage*, 90:423–435, 2014.
- [33] J. Canny. *The complexity of robot motion planning*. MIT press, 1988.
- [34] J. F. Canny. *The Complexity of Robot Motion Planning*. MIT Press, Cambridge, MA, USA, 1988.
- [35] J. F. Canny, E. Kaltofen, and L. Yagati. Solving systems of nonlinear polynomial equations faster. In *Proceedings of the ACM-SIGSAM 1989 International Symposium on Symbolic and Algebraic Computation*, ISSAC '89, page 121–128, New York, NY, USA, 1989. Association for Computing Machinery.
- [36] J. F. Canny, E. Kaltofen, and L. Yagati. Solving systems of nonlinear polynomial equations faster. In *Proceedings of the ACM-SIGSAM 1989 International Symposium on Symbolic and Algebraic Computation*, ISSAC '89, page 121–128, New York, NY, USA, 1989. Association for Computing Machinery.
- [37] J. Capco, M. Safey El Din, and J. Schicho. Robots, computer algebra and eight connected components. In *Proceedings of the 45th International Symposium on Symbolic and Algebraic Computation*, ISSAC '20, page 62–69, New York, NY, USA, 2020. Association for Computing Machinery.
- [38] C. Chauvin, M. Muller, and A. Weber. An application of quantifier elimination to mathematical biology. *Computer Algebra in Science and Engineering*. World Scientific, pages 287–296, 1994.
- [39] D. Chillingworth and M. Demazure. *Bifurcations and Catastrophes: Geometry of Solutions to Nonlinear Problems*. Universitext. Springer Berlin Heidelberg, 2013.

- [40] A. Chistov and D. Y. Grigor'ev. Polynomial-time factoring of multivariable polynomials over a general field. Technical report, Technical report, USSR Academy of Sciences, Steklov Mathematical Institute, 1982.
- [41] G. E. Collins. Quantifier elimination for real closed fields by Cylindrical Algebraic Decomposition: a synopsis. *ACM SIGSAM Bulletin*, 10(1):10–12, 1976.
- [42] G. E. Collins and H. Hong. Partial Cylindrical Algebraic Decomposition for quantifier elimination. *Journal of Symbolic Computation*, 12(3):299 – 328, 1991.
- [43] A. Conca and J. Herzog. On the Hilbert function of determinantal rings and their canonical module. *Proceedings of the American Mathematical Society*, 122(3):677–681, 1994.
- [44] R. Connelly and W. Whiteley. The stability of tensegrity frameworks. *International Journal of Space Structures*, 7(2):153–163, 1992.
- [45] S. Corvez and F. Rouillier. Using computer algebra tools to classify serial manipulators. In *International Workshop on Automated Deduction in Geometry*, pages 31–43. Springer, 09 2002.
- [46] O. Coss, J. Hauenstein, H. Hong, and D. Molzahn. Locating and counting equilibria of the kuramoto model with rank-one coupling. *SIAM Journal on Applied Algebra and Geometry*, 2, 04 2017.
- [47] M. Coste and M. Shiota. Thom's first isotopy lemma: a semialgebraic version, with uniform bounds. In F. Broglia, M. Galbiati, and A. Tognoli, editors, *Real Analytic and Algebraic Geometry*, page 83–101. De Gruyter, 1992.
- [48] D. A. Cox, J. Little, and D. O'Shea. *Ideals, Varieties, and Algorithms: An Introduction to Computational Algebraic Geometry and Commutative Algebra, 3/e (Undergraduate Texts in Mathematics)*. Springer-Verlag, Berlin, Heidelberg, 2007.
- [49] X. Dahan and É. Schost. Sharp estimates for triangular sets. In J. Gutierrez, editor, *Symbolic and Algebraic Computation, International Symposium ISSAC 2004, Santander, Spain, July 4-7, 2004, Proceedings*, pages 103–110. ACM, 2004.
- [50] J. H. Davenport and J. Heintz. Real quantifier elimination is doubly exponential. *J. Symb. Comput.*, 5(1):29 – 35, 1988.
- [51] L. De Moura and N. Bjørner. Z3: An efficient SMT solver. In C. R. Ramakrishnan and J. Rehof, editors, *Tools and Algorithms for the Construction and Analysis of Systems*, pages 337–340, Berlin, Heidelberg, 2008. Springer Berlin Heidelberg.
- [52] A. Dolzmann, T. Sturm, and V. Weispfenning. A new approach for automatic theorem proving in real geometry. *J. Autom. Reason.*, 21(3):357–380, Dec. 1998.

- [53] A. Dolzmann and V. Weispfenning. Local quantifier elimination. In *Proceedings of the 2000 International Symposium on Symbolic and Algebraic Computation*, page 86–94, New York, NY, USA, 2000. Association for Computing Machinery.
- [54] A. Dolzmann and L. A. Gilch. Generic hermitian quantifier elimination. In B. Buchberger and J. Campbell, editors, *Artificial Intelligence and Symbolic Computation*, pages 80–93, Berlin, Heidelberg, 2004. Springer Berlin Heidelberg.
- [55] D. Duval. Rational puiseux expansions. *Compositio Mathematica*, 70(2):119–154, 1989.
- [56] D. Eisenbud. *Commutative Algebra: With a View Toward Algebraic Geometry*. Graduate Texts in Mathematics. Springer, 1995.
- [57] J. Elliott, M. Giesbrecht, and É. Schost. On the bit complexity of finding points in connected components of a smooth real hypersurface. In I. Z. Emiris and L. Zhi, editors, *ISSAC '20: International Symposium on Symbolic and Algebraic Computation, Kalamata, Greece, July 20–23, 2020*, pages 170–177. ACM, 2020.
- [58] J. Faugère, P. Gaudry, L. Huot, and G. Renault. Polynomial systems solving by fast linear algebra. *CoRR*, abs/1304.6039, 2013.
- [59] J.-C. Faugère. A new efficient algorithm for computing Gröbner bases (F4). *Journal of pure and applied algebra*, 139(1-3):61–88, 1999.
- [60] J. C. Faugère. A new efficient algorithm for computing Gröbner bases without reduction to zero (F5). In *Proceedings of the 2002 international symposium on Symbolic and algebraic computation*, pages 75–83, 2002.
- [61] J.-C. Faugère, M. S. El Din, and P.-J. Spaenlehauer. Critical points and gröbner bases: The unmixed case. In *Proceedings of the 37th International Symposium on Symbolic and Algebraic Computation*, ISSAC '12, page 162–169, New York, NY, USA, 2012. Association for Computing Machinery.
- [62] J.-C. Faugère and A. Joux. Algebraic cryptanalysis of hidden field equation (HFE) cryptosystems using gröbner bases. In D. Boneh, editor, *Advances in Cryptology - CRYPTO 2003*, pages 44–60, Berlin, Heidelberg, 2003. Springer Berlin Heidelberg.
- [63] J.-C. Faugère, G. Moroz, F. Rouillier, and M. Safey El Din. Classification of the perspective-three-point problem, discriminant variety and real solving polynomial systems of inequalities. In *Proceedings of the Twenty-First International Symposium on Symbolic and Algebraic Computation*, ISSAC '08, page 79–86, New York, NY, USA, 2008. Association for Computing Machinery.

- [64] J.-C. Faugère, G. Moroz, F. Rouillier, and M. Safey El Din. Classification of the perspective-three-point problem, discriminant variety and real solving polynomial systems of inequalities. In *Proceedings of the twenty-first international symposium on Symbolic and algebraic computation*, pages 79–86, 2008.
- [65] J.-C. Faugère, M. Safey El Din, and T. Verron. On the complexity of computing Gröbner bases for quasi-homogeneous systems. In *Proceedings of the 38th International Symposium on Symbolic and Algebraic Computation*, ISSAC '13, page 189–196, New York, NY, USA, 2013. Association for Computing Machinery.
- [66] J.-C. Faugère. FGb: A Library for Computing Gröbner Bases. In K. Fukuda, J. Hoeven, M. Joswig, and N. Takayama, editors, *Mathematical Software - ICMS 2010*, volume 6327 of *Lecture Notes in Computer Science*, pages 84–87, Berlin, Heidelberg, September 2010. Springer Berlin / Heidelberg.
- [67] J.-C. Faugère, P. Gianni, D. Lazard, and T. Mora. Efficient computation of zero-dimensional gröbner bases by change of ordering. *Journal of Symbolic Computation*, 16(4):329–344, 1993.
- [68] J.-C. Faugère and C. Mou. Sparse FGLM algorithms. *Journal of Symbolic Computation*, 80:538–569, 2017.
- [69] J.-C. Faugère, M. Safey El Din, and P.-J. Spaenlehauer. On the complexity of the generalized Minrank problem. *Journal of Symbolic Computation*, 55:30 – 58, 2013.
- [70] J.-C. Faugère, M. Safey El Din, and T. Verron. On the complexity of computing Gröbner bases for weighted homogeneous systems. *Journal of Symbolic Computation*, 76:107–141, 2016.
- [71] I. A. Fotiou, P. Rostalski, P. A. Parrilo, and M. Morari. Parametric optimization and optimal control using algebraic geometry methods. *International Journal of Control*, 79(11):1340–1358, 2006.
- [72] R. Fukasaku, H. Iwane, and Y. Sato. Cgsqe/synrac: A real quantifier elimination package based on the computation of comprehensive gröbner systems. *ACM Commun. Comput. Algebra*, 50(3):101–104, Nov. 2016.
- [73] S. Galeani, D. Henrion, A. Jacquemard, and L. Zaccarian. Design of Marx generators as a structured eigenvalue assignment. *Automatica*, 50(10):2709–2717, 2014.
- [74] X.-S. Gao, X.-R. Hou, J. Tang, and H.-F. Cheng. Complete solution classification for the perspective-three-point problem. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(8):930–943, Aug. 2003.
- [75] N. Gaspar, X. Ren, C. Smith, J. Grima, and K. Evans. Novel honeycombs with auxetic behaviour. *Acta Materialia*, 53(8):2439 – 2445, 2005.

- [76] J. Gerhard, D. J. Jeffrey, and G. Moroz. A package for solving parametric polynomial systems. *ACM Commun. Comput. Algebra*, 43(3/4):61–72, June 2010.
- [77] É. Ghys and A. Ranicki. Signatures in algebra, topology and dynamics. *Ensaos Matemáticos*, 30:1 – 173, 2016.
- [78] P. M. Gianni and T. Teo Mora. Algebraic solution of systems of polynomial equations using Gröebner bases. In *Applied Algebra, Algebraic Algorithms and Error-Correcting Codes, 5th International Conference, AAECC-5, Menorca, Spain, June 15-19, 1987, Proceedings*, pages 247–257, 1987.
- [79] A. Giovini, T. Mora, G. Niesi, L. Robbiano, and C. Traverso. “one sugar cube, please” or selection strategies in the buchberger algorithm. In *Proceedings of the 1991 International Symposium on Symbolic and Algebraic Computation, ISSAC '91*, page 49–54, New York, NY, USA, 1991. Association for Computing Machinery.
- [80] M. Giusti. Some effectivity problems in polynomial ideal theory. In J. Fitch, editor, *EUROSAM 84*, pages 159–171, Berlin, Heidelberg, 1984. Springer Berlin Heidelberg.
- [81] M. Giusti. A note on the complexity of constructing standard bases. In B. F. Caviness, editor, *EUROCAL '85*, pages 411–412, Berlin, Heidelberg, 1985. Springer Berlin Heidelberg.
- [82] M. Giusti and J. Heintz. La détermination des points isolés et de la dimension d’une variété algébrique peut se faire en temps polynomial. *Computational Algebraic Geometry and Commutative Algebra*, 34:216–256, 02 1993.
- [83] M. Giusti, J. Heintz, J. Enrique Morais, and L. M. Pardo. Le rôle des structures de données dans les problèmes d’élimination. *Comptes Rendus de l’Académie des Sciences - Series I - Mathematics*, 325(11):1223–1228, 1997.
- [84] M. Giusti, J. Heintz, K. Hägele, J. E. Morais, L. M. Pardo, and J. L. Montaña. Lower bounds for diophantine approximation. In *In Proceedings of MEGA'96*, pages 277–317, 1997.
- [85] M. Giusti, J. Heintz, J. Morais, J. Morgenstem, and L. Pardo. Straight-line programs in geometric elimination theory. *Journal of Pure and Applied Algebra*, 124(1):101–146, 1998.
- [86] M. Giusti, J. Heintz, J. E. Morais, and L. M. Pardo. When polynomial equation systems can be “solved” fast? In *Applied Algebra, Algebraic Algorithms and Error-Correcting Codes, 11th International Symposium, AAECC-11, Paris, France, July 17-22, 1995, Proceedings*, pages 205–231, 1995.
- [87] M. Giusti, G. Lecerf, and B. Salvy. A Gröbner free alternative for polynomial system solving. *Journal of Complexity*, 17(1):154 – 211, 2001.

- [88] L. González-Vega, T. Recio, H. Lombardi, and M.-F. Roy. Sturm–Habicht sequences, determinants and real roots of univariate polynomials. In B. F. Caviness and J. R. Johnson, editors, *Quantifier Elimination and Cylindrical Algebraic Decomposition*, pages 300–316, Vienna, 1998. Springer Vienna.
- [89] D. R. Grayson and M. E. Stillman. Macaulay2, a software system for research in algebraic geometry. Available at <http://www.math.uiuc.edu/Macaulay2/>.
- [90] O. Grellier, P. Comon, B. Mourrain, and P. Trebuchet. Analytical blind channel identification. *IEEE Transactions on Signal Processing*, 50(9):2196–2207, 2002.
- [91] A. Greuet and M. Safey El Din. Probabilistic algorithm for polynomial optimization over a real algebraic set. *SIAM Journal on Optimization*, 24(3), 2014.
- [92] P. Griffiths and J. Harris. *Principles of algebraic geometry*. Wiley Classics Library. John Wiley & Sons, Inc., New York, 1994. Reprint of the 1978 original.
- [93] D. Y. Grigor’ev. Complexity of deciding Tarski algebra. *J. Symb. Comput.*, 5(1–2):65–108, Feb. 1988.
- [94] D. Y. Grigor’ev and N. Vorobjov. Solving systems of polynomial inequalities in subexponential time. *Journal of Symbolic Computation*, 5:37–64, 1988.
- [95] J. N. Grima and K. E. Evans. Auxetic behavior from rotating triangles. *Journal of Materials Science*, 41(10):3193–3196, May 2006.
- [96] L. Grégoire. Kronecker, 2002.
- [97] R. M. Hardt. Semi-algebraic local-triviality in semi-algebraic mappings. *American Journal of Mathematics*, 102(2):291–302, 1980.
- [98] A. Harnack. Ueber die Vieltheiligkeit der ebenen algebraischen Curven. *Math. Ann.*, 10(2):189–198, 1876.
- [99] K. Harris, J. D. Hauenstein, and A. Szanto. Smooth points on semi-algebraic sets, 2020.
- [100] R. Hartshorne. *Algebraic Geometry*. Graduate Texts in Mathematics. Springer New York, 2013.
- [101] J. Heintz. Definability and fast quantifier elimination in algebraically closed fields. *Theor. Comput. Sci.*, 24:239–277, 1983.
- [102] J. Heintz, M.-F. Roy, and P. Solernó. Sur la complexité du principe de tarski-seidenberg. *Bulletin de la Société Mathématique de France*, 118(1):101–126, 1990.
- [103] D. Henrion. Detecting rigid convexity of bivariate polynomials. *Linear Algebra and its Applications*, 432(5):1218 – 1233, 2010.

- [104] D. Henrion, S. Naldi, and M. Safey El Din. Real root finding for rank defects in linear hankel matrices. In K. Yokoyama, S. Linton, and D. Robertz, editors, *Proceedings of the 2015 ACM on International Symposium on Symbolic and Algebraic Computation, ISSAC 2015, Bath, United Kingdom, July 06 - 09, 2015*, pages 221–228. ACM, 2015.
- [105] D. Henrion and M. Sebek. Plane geometry and convexity of polynomial stability regions. In *Proceedings of the Twenty-First International Symposium on Symbolic and Algebraic Computation, ISSAC '08*, page 111–116, New York, NY, USA, 2008. Association for Computing Machinery.
- [106] C. Hermite. Sur le nombre des racines d’une équation algébrique comprises entre des limites données. extrait d’une lettre à m. borchardt. *J. Reine Angew. Math.*, 52:39–51, 1856.
- [107] H. Hong. An improvement of the projection operator in Cylindrical Algebraic Decomposition. In *Proceedings of the International Symposium on Symbolic and Algebraic Computation, ISSAC '90*, page 261–264, New York, NY, USA, 1990. Association for Computing Machinery.
- [108] H. Hong, R. Liska, and S. Steinberg. Testing stability by quantifier elimination. *Journal of Symbolic Computation*, 24(2):161–187, Aug. 1997.
- [109] H. Hong and M. Safey El Din. Variant real quantifier elimination: Algorithm and application. In *Proceedings of the 2009 International Symposium on Symbolic and Algebraic Computation, ISSAC '09*, page 183–190, New York, NY, USA, 2009. Association for Computing Machinery.
- [110] H. Hong and M. Safey El Din. Variant quantifier elimination. *J. Symb. Comput.*, 47(7):883 – 901, 2012. International Symposium on Symbolic and Algebraic Computation (ISSAC 2009).
- [111] J. Huisman. On the geometry of algebraic curves having many real components. *Rev. Mat. Complut.*, 14(1):83–92, 2001.
- [112] J. Huisman. On the geometry of algebraic curves having many real components. *Rev. Mat. Complut.*, 14(1):83–92, 2001.
- [113] J. Huisman. Non-special divisors on real algebraic curves and embeddings into real projective spaces. *Ann. Mat. Pura Appl. (4)*, 182(1):21–35, 2003.
- [114] C. G. Jacobi. Über eine elementare transformation eins in bezug auf jedes von zwei variablen-systemen linearen und homogenen ausdrucks. *Journal fur die reine und angewandte Mathematik* 53., pages 265 – 270, 1857.
- [115] Z. Jelonek. Testing sets for properness of polynomial mappings. *Mathematische Annalen*, 315:1–35, Sept. 1999.

- [116] Z. Jelonek and K. Kurdyka. Quantitative generalized bertini-sard theorem for smooth affine varieties. *Discrete & Computational Geometry*, 34(4):659–678, 2005.
- [117] J. Joos Heintz and J. Morgenstern. On the intrinsic complexity of elimination theory. *Journal of Complexity*, 9(4):471–498, 1993.
- [118] M. Kalkbrener. On the stability of gröbner bases under specializations. *Journal of Symbolic Computation*, 24(1):51–58, 1997.
- [119] G. Kanagalingam and N. Bajcinca. Mapping of eigenvalue performance specifications by real root classification. In *2018 26th Mediterranean Conference on Control and Automation (MED)*, pages 1–9, 2018.
- [120] D. Kapur. Automatically generating loop invariants using quantifier elimination. In *Dagstuhl Seminar Proceedings*. Schloss Dagstuhl-Leibniz-Zentrum für Informatik, 2006.
- [121] A. Kipnis and A. Shamir. Cryptanalysis of the HFE public key cryptosystem by relinearization. In M. Wiener, editor, *Advances in Cryptology - CRYPTO' 99*, pages 19–30, Berlin, Heidelberg, 1999. Springer Berlin Heidelberg.
- [122] A. Kobel, F. Rouillier, and M. Sagraloff. Computing real roots of real polynomials ... and now for real! In *Proceedings of the ACM on International Symposium on Symbolic and Algebraic Computation, ISSAC '16*, page 303–310, New York, NY, USA, 2016. Association for Computing Machinery.
- [123] V. A. Krasnov. Albanese mapping for *GMZ*-varieties. *Mat. Zametki*, 35(5):739–747, 1984.
- [124] L. Kronecker. Grundzüge einer arithmetischen theorie der algebraischen grössen. *Journal für die reine und angewandte Mathematik*, 92:1–122, 1882.
- [125] M. Kummer and D. Manevich. On Huisman’s conjectures about unramified real curves. *Preprint arXiv:1909.09601*, 2019.
- [126] Y. Kuramoto. Self-entrainment of a population of coupled non-linear oscillators. In H. Araki, editor, *International Symposium on Mathematical Problems in Theoretical Physics*, pages 420–422, Berlin, Heidelberg, 1975. Springer Berlin Heidelberg.
- [127] K. Kurdyka, P. Orro, and S. Simon. Semialgebraic sard theorem for generalized critical values. *Journal of Differential Geometry*, 56:67–92, 09 2000.
- [128] G. Lafferriere, G. Pappas, G. Schneider, and S. Yovine. Parameter synthesis in robot motion planning using symbolic reachability computations. In *Proceedings of the 8th IEEE Mediterranean Conference on Control and Automation*. Citeseer, 2000.
- [129] R. Lakes. Foam structures with a negative poisson’s ratio. *Science*, 235(4792):1038–1040, 1987.

- [130] Y. N. Lakshman. *A Single Exponential Bound on the Complexity of Computing Gröbner Bases of Zero Dimensional Ideals*, pages 227–234. Birkhäuser Boston, Boston, MA, 1991.
- [131] Y. N. Lakshman and D. Lazard. On the complexity of zero-dimensional algebraic systems. In *Effective methods in algebraic geometry*, pages 217–225. Springer, 1991.
- [132] J. B. Lasserre. Global optimization with polynomials and the problem of moments. *SIAM Journal on Optimization*, 11:796–817, 2001.
- [133] D. Lazard. Gröbner bases, gaussian elimination and resolution of systems of algebraic equations. In *European Conference on Computer Algebra*, pages 146–156. Springer, 1983.
- [134] D. Lazard and F. Rouillier. Solving parametric polynomial systems. *Journal of Symbolic Computation*, 42(6):636–667, 2007.
- [135] H. P. Le, D. Manevich, and D. Plaumann. Computing totally real hyperplane sections and linear series on algebraic curves. *Le Matematiche*, 2021.
- [136] H. P. Le and M. Safey El Din. Solving parametric systems of polynomial equations over the reals through Hermite matrices. Preprint in the review process, Nov. 2020.
- [137] H. P. Le and M. Safey El Din. Faster One Block Quantifier Elimination for Regular Polynomial Systems of Equations. In *International Symposium on Symbolic and Algebraic Computation 2021 (ISSAC '21)*, Saint Petersburg, Russia, July 2021.
- [138] H. P. Le, M. Safey El Din, and T. de Wolff. Computing the real isolated points of an algebraic hypersurface. In I. Z. Emiris and L. Zhi, editors, *ISSAC '20: International Symposium on Symbolic and Algebraic Computation, Kalamata, Greece, July 20-23, 2020*, pages 297–304. ACM, 2020.
- [139] G. Lecerf. Computing an equidimensional decomposition of an algebraic variety by means of geometric resolutions. In *Proceedings of the 2000 International Symposium on Symbolic and Algebraic Computation, ISSAC 2000, St. Andrews, United Kingdom, August 6-10, 2000*, pages 209–216, 2000.
- [140] T.-Y. Li. Numerical solution of multivariate polynomial systems by homotopy continuation methods. *Acta Numerica*, 6:399–436, 1997.
- [141] T.-Y. Li and X. Wang. Solving real polynomial systems with real homotopies. *Mathematics of Computation*, 60(202):669–680, 1993.
- [142] S. Liang and D. J. Jeffrey. Automatic computation of the complete root classification for a parametric polynomial. *Journal of Symbolic Computation*, 44(10):1487–1501, 2009.

- [143] S. Liang, D. J. Jeffrey, and M. M. Maza. The complete root classification of a parametric polynomial on an interval. In *Proceedings of the twenty-first international symposium on Symbolic and algebraic computation*, pages 189–196, 2008.
- [144] R. Liska and S. L. Steinberg. Applying Quantifier Elimination to Stability Analysis of Difference Schemes. *The Computer Journal*, 36(5):497–503, 01 1993.
- [145] Q. Liu. *Algebraic geometry and arithmetic curves*, volume 6 of *Oxford Graduate Texts in Mathematics*. Oxford University Press, Oxford, 2002. Translated from the French by Reinie Ern e, Oxford Science Publications.
- [146] E. W. Mayr and A. R. Meyer. The complexity of the word problems for commutative semi-groups and polynomial ideals. *Advances in Mathematics*, 46(3):305–329, 1982.
- [147] M. M. Maza, B. Xia, and R. Xiao. On solving parametric polynomial systems. *Mathematics in Computer Science*, 6(4):457–473, 2012.
- [148] S. McCallum. An improved projection operation for Cylindrical Algebraic Decomposition of three-dimensional space. *Journal of Symbolic Computation*, 5(1):141 – 161, 1988.
- [149] S. McCallum. On projection in CAD-based quantifier elimination with equational constraint. In *Proceedings of the 1999 International Symposium on Symbolic and Algebraic Computation*, ISSAC ’99, page 145–149, New York, NY, USA, 1999. Association for Computing Machinery.
- [150] J. Milnor. On the Betti numbers of real varieties. *Proceedings of the American Mathematical Society*, 15:275–280, 1964.
- [151] H. M. M oller and F. Mora. Upper and lower bounds for the degree of groebner bases. In J. Fitch, editor, *EUROSAM 84*, pages 172–183, Berlin, Heidelberg, 1984. Springer Berlin Heidelberg.
- [152] J.-P. Monnier. Divisors on real curves. *Adv. Geom.*, 3(3):339–360, 2003.
- [153] J.-P. Monnier. On real generalized Jacobian varieties. *J. Pure Appl. Algebra*, 203(1-3):252–274, 2005.
- [154] A. Montes and M. Wibmer. Gr obner bases for polynomial systems with parameters. *Journal of Symbolic Computation*, 45(12):1391 – 1425, 2010. MEGA’2009.
- [155] G. Moreno-Socias. Degrevlex gr obner bases of generic complete intersections. *Journal of Pure and Applied Algebra*, 180(3):263 – 283, 2003.
- [156] C. B. Mulligan, J. H. Davenport, and M. England. Theoryguru: A mathematica package to apply quantifier elimination technology to economics. In *International Congress on Mathematical Software*, pages 369–378. Springer, 2018.

- [157] W. Niu and D. Wang. Algebraic approaches to stability analysis of biological systems. *Mathematics in Computer Science*, 1(11):507–539, 2008.
- [158] L. M. Pardo. How lower and upper complexity bounds meet in elimination theory. In G. Cohen, M. Giusti, and T. Mora, editors, *Applied Algebra, Algebraic Algorithms and Error-Correcting Codes*, pages 33–69, Berlin, Heidelberg, 1995. Springer Berlin Heidelberg.
- [159] K. Pardue. Generic sequences of polynomials. *Journal of Algebra*, 324(4):579 – 590, 2010.
- [160] P. Pedersen, M.-F. Roy, and A. Szpirglas. Counting real zeros in the multivariate case. In F. Eyssette and A. Galligo, editors, *Computational Algebraic Geometry*, pages 203–224, Boston, MA, 1993. Birkhäuser Boston.
- [161] D. Plaumann, B. Sturmfels, and C. Vinzant. Quartic curves and their bitangents. *J. Symbolic Comput.*, 46(6):712–733, 2011.
- [162] A. Popolitov and S. Shakirov. On undulation invariants of plane curves. *Michigan Math. J.*, 64(1):143–153, 2015.
- [163] M. Raghavan and B. Roth. Solving Polynomial Systems for the Kinematic Analysis and Synthesis of Mechanisms and Robot Manipulators. *Journal of Mechanical Design*, 117(B):71–79, 06 1995.
- [164] J. Renegar. On the computational complexity and geometry of the first-order theory of the reals. Part III: Quantifier elimination. *J. Symb. Comput.*, 13(3):329–352, Mar. 1992.
- [165] B. Roth and W. Whiteley. Tensegrity frameworks. *Trans. Am. Math. Soc.*, 265:419–446, 1981.
- [166] F. Rouillier. Solving zero-dimensional systems through the rational univariate representation. *Applicable Algebra in Engineering, Communication and Computing*, 9(5):433–461, 1999.
- [167] F. Rouillier, M. Roy, and M. Safey El Din. Finding at least one point in each connected component of a real algebraic set defined by a single equation. *J. Complexity*, 16(4):716–750, 2000.
- [168] F. Rouillier and P. Zimmermann. Efficient isolation of polynomial’s real roots. *Journal of Computational and Applied Mathematics*, 162(1):33–50, Jan. 2004.
- [169] M. Safey El Din. Finding sampling points on real hypersurfaces is easier in singular situations. In *MEGA*, 2005.
- [170] M. Safey El Din. Real algebraic geometry library, RAGlib (version 3.4), 2017.
- [171] M. Safey El Din and E. Schost. Polar varieties and computation of one point in each connected component of a smooth real algebraic set. In *Proc. of the 2003 Int. Symp. on Symb. and Alg. Comp.*, ISSAC ’03, page 224–231. ACM, 2003.

- [172] M. Safey El Din and É. Schost. Properness defects of projections and computation of at least one point in each connected component of a real algebraic set. *Discrete & Computational Geometry*, 32(3):417–430, 2004.
- [173] M. Safey El Din and É. Schost. A baby steps/giant steps probabilistic algorithm for computing roadmaps in smooth bounded real hypersurface. *Discrete & Computational Geometry*, 45(1):181–220, 2011.
- [174] M. Safey El Din and É. Schost. A nearly optimal algorithm for deciding connectivity queries in smooth and bounded real algebraic sets. *J. ACM*, 63(6):48:1–48:37, Jan. 2017.
- [175] M. Safey El Din and E. Schost. Bit complexity for multi-homogeneous polynomial system solving—application to polynomial minimization. *Journal of Symbolic Computation*, 87:176–206, 2018.
- [176] M. Sagraloff. On the complexity of the Descartes method when using approximate arithmetic. *Journal of Symbolic Computation*, 65:79–110, 2014.
- [177] C. Scheiderer. Sums of squares of regular functions on real algebraic varieties. *Trans. Amer. Math. Soc.*, 352(3):1039–1069, 2000.
- [178] É. Schost. Computing parametric geometric resolutions. *Applicable Algebra in Engineering, Communication and Computing*, 13(5):349–393, 2003.
- [179] É. Schost. Computing parametric geometric resolutions. *Applicable Algebra in Engineering, Communication and Computing*, 13(5):349–393, 2003.
- [180] J. T. Schwartz and M. Sharir. On the “piano movers” problem. II. general techniques for computing topological properties of real algebraic manifolds. *Advances in Applied Mathematics*, 4(3):298 – 351, 1983.
- [181] K. Schwede and Z. Yang. Divisor package for Macaulay2. *J. Softw. Algebra Geom.*, 8:87–94, 2018.
- [182] A. Seidenberg. A new decision method for elementary algebra. *Annals of Mathematics*, 60(2):365–374, 1954.
- [183] A. Seidl and T. Sturm. A generic projection operator for partial cylindrical algebraic decomposition. In *Proceedings of the 2003 International Symposium on Symbolic and Algebraic Computation*, ISSAC ’03, page 240–247, New York, NY, USA, 2003. Association for Computing Machinery.
- [184] I. R. Shafarevich. *Basic Algebraic Geometry 1: Varieties in Projective Space*. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013.

- [185] G. I. Sivashinsky. Diffusional-thermal theory of cellular flames. *Combustion Science and Technology*, 15(3-4):137–145, 1977.
- [186] S. Soliman, F. Fages, and O. Radulescu. A constraint solving approach to model reduction by tropical equilibration. *Algorithms for Molecular Biology*, 9:24–34, 4 2014.
- [187] P.-J. Spaenlehauer. On the complexity of computing critical points with Gröbner bases. *SIAM Journal on Optimization*, 24:1382–1401, 07 2014.
- [188] C. Sturm. Mémoire sur la résolution des équations numériques. *Mémoires présentés à L’Institut des Sciences, Lettres et Arts, par divers savants et lus dans ses assemblées : Sciences, Mathématiques et Physiques*, pages 273–318, 1835.
- [189] T. Sturm and A. Tiwari. Verification and synthesis using real quantifier elimination. In *Proceedings of the 36th International Symposium on Symbolic and Algebraic Computation, ISSAC ’11*, page 329–336, New York, NY, USA, 2011. Association for Computing Machinery.
- [190] T. Sturm and V. Weispfenning. Computational geometry problems in REDLOG. In *Selected Papers from the International Workshop on Automated Deduction in Geometry*, page 58–86, Berlin, Heidelberg, 1996. Springer-Verlag.
- [191] J. J. Sylvester. A demonstration of the theorem that every homogeneous quadratic polynomial is reducible by real orthogonal substitution to the form of a sum of positive and negative squares. *Philosophical Magazine IV.*, pages 138 – 142, 1852.
- [192] A. Tarski. *A Decision Method for Elementary Algebra and Geometry*. University of California Press, 1951.
- [193] R. Thom. *Sur L’Homologie des Varietes Algebriques Réelles*, pages 255–265. Princeton University Press, 1965.
- [194] T. Verron. *Regularisation of Gröbner basis computations for weighted and determinantal systems, and application to medical imagery*. Theses, Université Pierre et Marie Curie - Paris VI, Sept. 2016.
- [195] J. Verschelde. Polynomial homotopy continuation with phcpack. *ACM Communications in Computer Algebra*, 44(3/4):217–220, Jan. 2011.
- [196] N. Vorobjov. Complexity of computing the local dimension of a semialgebraic set. *Journal of Symbolic Computation*, 27(6):565–579, 1999.
- [197] V. Weispfenning. The complexity of linear problems in fields. *Journal of Symbolic Computation*, 5(1):3–27, 1988.
- [198] V. Weispfenning. Comprehensive gröbner bases. *Journal of Symbolic Computation*, 14(1):1–29, 1992.

- [199] V. Weispfenning. A new approach to quantifier elimination for real algebra. In *Quantifier Elimination and Cylindrical Algebraic Decomposition*, pages 376–392, Vienna, 1998. Springer Vienna.
- [200] B. Xia and L. Yang. An algorithm for isolating the real solutions of semi-algebraic systems. *Journal of Symbolic Computation*, 34(5):461–477, 2002.
- [201] L. Yang. A simplified algorithm for solution classification of the perspective-three-point problem. In *Mathematics-Mechanization Research Preprints*, 12 1998.
- [202] L. Yang, X. Hou, and B. Xia. A complete algorithm for automated discovering of a class of inequality-type theorems. *Science in China Series F Information Sciences*, 44(1):33–49, 2001.
- [203] L. Yang and B. Xia. Real solution classification for parametric semi-algebraic systems. In A. Dolzmann, A. Seidl, and T. Sturm, editors, *Algorithmic Algebra and Logic. Proceedings of the A3L 2005, April 3-6, Passau, Germany; Conference in Honor of the 60th Birthday of Volker Weispfenning*, pages 281–289, 2005.
- [204] L. Yang and Z. Zeng. Equi-cevaline points on triangles. In *Computer Mathematics: Proceedings of the Fourth Asian Symposium (ASCM 2000)*, page 130. World Scientific Publishing Company Incorporated, 2000.
- [205] L. Yang and Z. Zeng. An open problem on metric invariants of tetrahedra. In *Proceedings of the 2005 International Symposium on Symbolic and Algebraic Computation, ISSAC '05*, page 362–364, New York, NY, USA, 2005. Association for Computing Machinery.
- [206] L. Yang, N. Zhan, B. Xia, and C. Zhou. *Program Verification by Using DISCOVERER*, pages 528–538. Springer Berlin Heidelberg, Berlin, Heidelberg, 2008.
- [207] W. Yang, Z.-M. Li, W. Shi, and B.-H. Xie. Review on auxetic materials. *Journal of Materials Science*, 39:3269–3279, 2004.