

The complexity and geometry of numerically solving polynomial systems.

Carlos Beltrán and Michael Shub

This paper is dedicated to the memory of our beloved friend and colleague Jean Pierre Dedieu.

ABSTRACT. These pages contain a short overview on the state of the art of efficient numerical analysis methods that solve systems of multivariate polynomial equations. We focus on the work of Steve Smale who initiated this research framework, and on the collaboration between Stephen Smale and Michael Shub, which set the foundations of this approach to polynomial system-solving, culminating in the more recent advances of Carlos Beltrán, Luis Miguel Pardo, Peter Bürgisser and Felipe Cucker.

1. The modern numerical approach to polynomial system solving

In this paper we survey some of the recent advances in the solution of polynomial systems. Such a classical topic has been studied by hundreds of authors from many different perspectives. We do not intend to make a complete historical description of all the advances achieved during the last century or two, but rather to describe in some detail the state of the art of what we think is the most successful (both from practical and theoretical perspectives) approach. Homotopy methods are used to solve polynomial systems in real life applications all around the world.

The key ingredient of homotopy methods is a one-line thought: given a goal system to be solved, choose some other system (similar in form, say with the same degree and number of variables) with a known solution ζ_0 , and move this new system to the goal system, tracking how the known solution moves to a solution of the goal. Before stating any notation, we can explain briefly why this process is reasonable: if for every $t \in [0, 1]$ we have a system of equations f_t (f_0 is the system with a known solution, f_1 is the one we want to solve), then we are looking for a path ζ_t , $t \in [0, 1]$, such that $f_t(\zeta_t) = 0$. As long as the derivative $df_t(\zeta_t)$ is invertible for all t we can continue the solution from f_0 to f_1 , by the implicit function theorem. Now we have various methods to accomplish this continuation. We can slowly increment t and use iterative numerical solution methods such as Newton's method to track the solution or we may differentiate the expression $f_t(\zeta_t) = 0$ and solve for

2010 *Mathematics Subject Classification.* Primary 65H10, 14Q20, 68Q25.

The first author was partially Supported by MTM2010-16051, Spanish Ministry of Science.

The second author was partially supported by a CONICET grant PIP0801 2010-2012 and by ANPCyT PICT 2010-00681.

We thank an anonymous referee for his detailed reading of the manuscript.

$d/(dt)(\zeta_t) = \dot{\zeta}_t$. Then, we can write our problem as an initial value problem:

$$(1.1) \quad \begin{cases} \dot{\zeta}_t = -Df_t(\zeta_t)^{-1}f_t(\zeta_t) \\ \zeta_0 \text{ known} \end{cases}$$

Systems of ODEs have been much studied and hence this is an interesting idea: we have reduced our original problem to a very much studied one. One can just plug in a standard numerical ODE solver such as backward Euler or a version of Runge–Kutta’s method. Even then, in practice, it is desirable to, from time to time, perform some steps of Newton’s method $z \rightarrow x - Df_t(x)^{-1}f_t(x)$ to our approximation z_t of ζ_t , to get closer to the path (f_t, ζ_t) . After some testing and adjustment of parameters, this naïve idea can be made to work with impressive practical performance and there are several software packages which attain spectacular results (solving systems with many variables and high degree) in a surprisingly short running time, see for example [7, 41, 42, 63]

From a mathematical point of view, there are several things in the process we have just described that need to be analyzed: will there actually exist a path ζ_t (maybe it is only defined for, say, $t < 1/2$)? what is the expected complexity of the process (in particular, can we expect average polynomial running time in some sense)? what “simple system with a known solution” should we start at? how should we join f_0 and f_1 , that is what should be the path f_t ?

In the last few decades a lot of progress has been made in studying these questions. This progress is the topic of this paper.

2. A technical description of the problem

We will center our attention in Smale’s 17–th problem, which we recall now.

PROBLEM 2.1. *Can a zero of n complex polynomial equations in n unknowns be found approximately, on the average, in polynomial time with a uniform algorithm?*

We have written in bold the technical terms that need to be clarified.

In order to understand the details of the problem and the solution suggested in Section 1, we need to describe some important concepts and notation in detail. Maybe the first one is our understanding of what a “solution” is: clearly, one cannot expect solutions of polynomial systems to be rational numbers, so one can only search for “quasi–solutions” in some sense. There are several definitions of such a thing, the most stable being the following one (introduced in [57], see also [23, 39, 40]):

DEFINITION 2.2. Given a polynomial system, understood as a mapping $f : \mathbb{C}^n \rightarrow \mathbb{C}^n$, an approximate zero of f with associated (exact) zero ζ is a vector $z_0 \in \mathbb{C}^n$ such that

$$\|z_k - \zeta\| \leq \frac{1}{2^{2^k-1}}\|z_0 - \zeta\|, \quad k \geq 0,$$

where z_k is the result of applying k times Newton’s operator $z \mapsto z - Df(z)^{-1}f(z)$ (note that the definition of approximate zero implicitly assumes that z_k is defined for all $k \geq 0$.)

The power of this definition is that, as we will see below, given any polynomial system f and any exact zero $\zeta \in \mathbb{C}^n$, approximate zeros of f with associated zero ζ exist whenever $Df(\zeta)$ is an invertible matrix.

Recall that our first goal is to transform the problem of polynomial system solving into an implicit function problem or an ODE system like that of (1.1). There exist two principal reasons why the solution of such a system can fail to be defined for all $t > 0$: that the function defining the derivative is not everywhere defined (this corresponds naturally to $Df_t(\zeta_t)$ not being invertible), and that the solution escapes to infinity. The first problem seems to be more delicate and difficult to solve, but the second one is actually very easily dealt with: we just need to define our ODE in a compact manifold, instead of just in \mathbb{C}^n . The most similar compact manifold to \mathbb{C}^n is $\mathbb{P}(\mathbb{C}^{n+1})$, and the way to take the problem into $\mathbb{P}(\mathbb{C}^{n+1})$ is just homogenizing the equations.

DEFINITION 2.3. Let $f : \mathbb{C}^n \rightarrow \mathbb{C}^n$ be a polynomial system, that is $f = (f_1, \dots, f_n)$ where $f_i : \mathbb{C}^n \rightarrow \mathbb{C}$ is a polynomial of degree some d_i ,

$$f(x_1, \dots, x_n) = \sum_{\alpha_1 + \dots + \alpha_n \leq d_i} a_{\alpha_1, \dots, \alpha_n}^{(i)} x_1^{\alpha_1} \dots x_n^{\alpha_n}.$$

The homogeneous counterpart of f is $h : \mathbb{C}^{n+1} \rightarrow \mathbb{C}^n$ defined by $h = (h_1, \dots, h_n)$ where

$$h(x_0, x_1, \dots, x_n) = \sum_{\alpha_1 + \dots + \alpha_n \leq d_i} a_{\alpha_1, \dots, \alpha_n}^{(i)} x_0^{d_i - \sum_{i=1}^n \alpha_i} x_1^{\alpha_1} \dots x_n^{\alpha_n}.$$

We will talk about such a system h simply as a homogeneous system.

Note that if ζ is a zero of f then $(1, \zeta)$ is a zero of the homogeneous counterpart h of f . Reciprocally, if $\zeta = (\zeta_0, \zeta_1, \dots, \zeta_n)$ is a zero of h and if $\zeta_0 \neq 0$, then $(\zeta_1/\zeta_0, \dots, \zeta_n/\zeta_0)$ is a zero of f . Thus, the zeros of f and h are in correspondence and we can think of solving h and then recovering the zeros of f (this is not a completely obvious process when we only have approximate zeros, see [15].) Moreover, it is clear that for any complex number $\lambda \in \mathbb{C}$ and for $x \in \mathbb{C}^{n+1}$ we have

$$h(\lambda x) = \text{Diag}(\lambda^{d_1}, \dots, \lambda^{d_n})h(x),$$

and thus the zeros of h lie naturally in the projective space $\mathbb{P}(\mathbb{C}^{n+1})$.

As we will be working with homogeneous systems and projective zeros, we need a definition of approximate zero in the spirit of Definition 2.2 which is amenable to a projective setting. The following one, which uses the projective version [50] of Newton's operator, makes the work. Here and throughout the paper, given a matrix or vector A , by A^* we mean the complex conjugate transpose of A , and by $d_R(x, y)$ we mean the Riemannian distance from x to y , where x and y are elements in some Riemannian manifold.

DEFINITION 2.4. Given a homogeneous system h , an approximate zero of h with associated (exact) zero $\zeta \in \mathbb{P}(\mathbb{C}^{n+1})$ is a vector $z_0 \in \mathbb{P}(\mathbb{C}^{n+1})$ such that

$$d_R(z_k, \zeta) \leq \frac{1}{2^{k-1}} d_R(z_0, \zeta), \quad k \geq 0,$$

where z_k is the result of applying k times the projective Newton operator $z \mapsto z - Dh(z)|_{z^\perp}^{-1} h(z)$ (again, the definition of approximate zero implicitly assumes that z_k is defined for all $k \geq 0$.) Here, by $Dh(z)|_{z^\perp}$ we mean the restriction of the derivative of h at z , to the (complex) orthogonal subspace $z^\perp = \{y \in \mathbb{C}^{n+1} : y^* z = 0\}$.

It is a simple exercise to verify that (projective) Newton's method is well defined, that is the point it defines in projective space does not depend on the representative $z \in \mathbb{C}^{n+1}$ chosen for a point in projective space.

A (projective) approximate zero of h is thus a projective point such that the successive iterates of the projective Newton operator quickly approach an exact zero of h . Thus finding an approximate zero is an excellent output of a numerical zero-finding algorithm to solve h .

Because we are going to consider paths of systems $\{h_t\}_{t \in [a,b]}$, it is convenient to fix a framework where one can define these nicely. To this end, we consider the vector space of homogeneous polynomials of fixed degree $s \geq 1$:

$$\mathcal{H}_s = \{h \in \mathbb{C}[x_0, \dots, x_n] : h \text{ is homogeneous of degree } s\}.$$

It is convenient to consider an Hermitian product (and the associated metric) on \mathcal{H}_s . A desirable property of such a metric is the unitary invariance, namely, we would like to have an Hermitian product such that

$$\langle h, g \rangle_{\mathcal{H}_s} = \langle h \circ U, g \circ U \rangle_{\mathcal{H}_s}, \quad \forall U \in \mathcal{U}_{n+1},$$

where \mathcal{U}_{n+1} is the group of unitary matrices of size $n+1$. Such property was studied in detail in [52]. It turns out that there exists a unique (up to scalar multiplication) Hermitian product that satisfies it, the one defined as follows:

$$\begin{aligned} \left\langle \sum_{\alpha_0 + \dots + \alpha_n = s} a_{\alpha_0, \dots, \alpha_n} x_0^{\alpha_0} \dots x_n^{\alpha_n}, \sum_{\alpha_0 + \dots + \alpha_n = s} b_{\alpha_0, \dots, \alpha_n} x_0^{\alpha_0} \dots x_n^{\alpha_n} \right\rangle_{\mathcal{H}_s} = \\ \sum_{\alpha_0 + \dots + \alpha_n = s} \frac{\alpha_0! \dots \alpha_n!}{s!} a_{\alpha_0, \dots, \alpha_n} \overline{b_{\alpha_0, \dots, \alpha_n}}, \end{aligned}$$

where $\bar{\cdot}$ just means complex conjugation. Note that this is just a weighted version of the standard complex Hermitian product in complex affine space.

Then, given a list of degrees $(d) = (d_1, \dots, d_n)$, we consider the vector space

$$\mathcal{H}_{(d)} = \prod_{i=1}^n \mathcal{H}_{d_i}.$$

Note that an element h of $\mathcal{H}_{(d)}$ can be seen both as a mapping $h : \mathbb{C}^{n+1} \rightarrow \mathbb{C}^n$ or as a polynomial system, and can be identified by the list of coefficients of h_1, \dots, h_n . We denote by $\mathbb{P}(\mathcal{H}_{(d)})$ the projective space associated to $\mathcal{H}_{(d)}$, by N the complex dimension of $\mathbb{P}(\mathcal{H}_{(d)})$ (so the dimension of $\mathcal{H}_{(d)}$ is $N+1$) and we consider the following Hermitian structure in $\mathcal{H}_{(d)}$:

$$\langle h, g \rangle = \sum_{i=1}^n \langle h_i, g_i \rangle_{\mathcal{H}_{d_i}}, \quad \|h\| = \langle h, h \rangle^{1/2}.$$

This Hermitian product (and the associate Hermitian structure and metric) is also called the Bombieri–Weyl or the Kostlan product (structure, metric). As usual, this Hermitian product in $\mathcal{H}_{(d)}$ defines an associated Riemannian structure given by the real part of $\langle \cdot, \cdot \rangle$. We can thus consider integrals of functions defined on $\mathcal{H}_{(d)}$.

We denote by \mathbb{S} the unit sphere in $\mathcal{H}_{(d)}$, and we endow \mathbb{S} with the inherited Riemannian structure from that of $\mathcal{H}_{(d)}$. Then, $\mathbb{P}(\mathcal{H}_{(d)})$ has a natural Riemannian structure, the unique one making the projection $\mathbb{S} \rightarrow \mathbb{P}(\mathcal{H}_{(d)})$ a Riemannian submersion. That is the derivative of the projection restricted to the normal to the

fibers is an isometry. We can thus also consider integrals of functions defined in \mathbb{S} or $\mathbb{P}(\mathcal{H}_{(d)})$. We can now talk about probabilities in \mathbb{S} or $\mathbb{P}(\mathcal{H}_{(d)})$: given a measurable (nonnegative or integrable) mapping X defined in \mathbb{S} or $\mathbb{P}(\mathcal{H}_{(d)})$, we can consider its expected value:

$$E_{\mathbb{S}}(X) = \frac{1}{\nu(\mathbb{S})} \int_{\mathbb{S}} X(h) dh \quad \text{or} \quad E_{\mathbb{P}(\mathcal{H}_{(d)})}(X) = \frac{1}{\nu(\mathbb{P}(\mathcal{H}_{(d)}))} \int_{\mathbb{P}(\mathcal{H}_{(d)})} X(h) dh,$$

where we simply denote by $\nu(E)$ the volume of a Riemannian manifold E . Similarly, one can talk about probabilities in $\mathcal{H}_{(d)}$ according to the standard Gaussian distribution compatible with $\langle \cdot, \cdot \rangle$: given a measurable (nonnegative or integrable) mapping X defined in $\mathcal{H}_{(d)}$, its expected value is:

$$E_{\mathcal{H}_{(d)}}(X) = \frac{1}{(2\pi)^{N+1}} \int_{\mathcal{H}_{(d)}} X(h) e^{-\|h\|^2/2} dh.$$

We can now come back to Problem 2.1 and see what do each of the terms in that problem mean: Smale himself points out that one can just solve homogeneous systems (as suggested above). We still have a few terms to clarify:

- found approximately. This means finding an approximate zero in the sense of Definition 2.4.
- on the average, in polynomial time. This now means that, if $X(h)$ is the time needed by the algorithm to output an approximate zero of the input system h , then the expected value of X is a quantity polynomial in the input size, that is polynomial in N . The number of variables, n , and the maximum of the degrees, d , are smaller than N , and hence one attempts to get a bound on the expected value of X , as a polynomial in n, d, N .
- uniform algorithm. Smale demands an algorithm in the Blum–Shub–Smale model [20, 21], that is exact operations and comparisons between real numbers are assumed. This assumption departs from the actual performance of our computers, but it is close enough to be translated to performance in many situations. Uniform means that the same algorithm works for all (d) and n .

3. Geometry and condition number

We can now set up a geometric framework for homotopy methods. Consider the following set, usually called the solution variety:

$$(3.1) \quad \mathcal{V} = \{(h, \zeta) \in \mathbb{P}(\mathcal{H}_{(d)}) \times \mathbb{P}(\mathbb{C}^{n+1}) : h(\zeta) = 0\}.$$

This set is actually a smooth complex submanifold (as well as a complex algebraic subvariety) of $\mathbb{P}(\mathcal{H}_{(d)}) \times \mathbb{P}(\mathbb{C}^{n+1})$, see [20], and is clearly compact. It will be useful to consider the following diagram.

$$(3.2) \quad \begin{array}{ccc} & \mathcal{V} & \\ \pi_1 \swarrow & & \searrow \pi_2 \\ \mathbb{P}(\mathcal{H}_{(d)}) & & \mathbb{P}(\mathbb{C}^{n+1}) \end{array}$$

It is clear that $\pi_1^{-1}(h)$ is a copy of the zero set of h . Reciprocally, for fixed $\zeta \in \mathbb{P}(\mathbb{C}^{n+1})$, the set $\pi_2^{-1}(\zeta)$ is the vector space of polynomial systems that have ζ as a zero.

Let $\Sigma' \subseteq \mathcal{V}$ be the set of critical points of π_1 and $\Sigma = \pi_1(\Sigma') \subseteq \mathbb{P}(\mathcal{H}_{(d)})$ the set of critical values of π_1 . It is not hard to prove that:

- π_1 restricted to the set $\mathcal{V} \setminus \pi_1^{-1}(\Sigma)$ is a (smooth) \mathcal{D} -fold covering map, where $\mathcal{D} = d_1 \cdots d_n$ is the Bezóut number.
- $\Sigma' = \{(h, \zeta) \in \mathcal{V} : Dh(\zeta) |_{\zeta^\perp} \text{ has non-maximal rank}\}$. In that case, we say that ζ is a singular zero of h . Otherwise, we say that ζ is a regular zero of h .

This means, in particular, that the homotopy process described above can be carried out whenever the path of systems lies outside of Σ :

THEOREM 3.1. *Let $\{h_t : t \in [a, b]\}$ be a C^1 curve in $\mathbb{P}(\mathcal{H}_{(d)}) \setminus \Sigma$ and let ζ be a zero of h_a . Then, there exists a unique lift of h_t through π_1 , that is a C^1 curve $(h_t, \zeta_t) \in \mathcal{V}$ such that $\zeta_a = \zeta$. In particular, ζ_b is a zero of h_b . Moreover, the lifted curve satisfies:*

$$(3.3) \quad \frac{d}{dt}(h_t, \zeta_t) = \left(\dot{h}_t, -Dh_t(\zeta_t) |_{\zeta_t^\perp}^{-1} \dot{h}_t(\zeta_t) \right).$$

Finally, the set $\Sigma \subseteq \mathbb{P}(\mathcal{H}_{(d)})$ is a complex projective algebraic variety, thus it has real codimension 2 and the projection of most real lines in $\mathcal{H}_{(d)}$ to $\mathbb{P}(\mathcal{H}_{(d)})$ does not intersect Σ .

The last claim of Theorem 3.1 must be understood as follows. Let $g, f \in \mathcal{H}_{(d)}$ be chosen at random. Then, with probability one, the projection to $\mathbb{P}(\mathcal{H}_{(d)})$ of the line containing g and f does not intersect Σ .

In the case the thesis of Theorem 3.1 holds we just say that ζ_a can be continued to a zero ζ_b of h_a . One can be even more precise:

THEOREM 3.2. *Let $\{h_t : t \in [a, b]\}$ be a C^1 curve in $\mathbb{P}(\mathcal{H}_{(d)}) \setminus \Sigma$ and let ζ be a zero of h_a . Then, every zero ζ of h_a can be continued to a zero of h_b , defining a bijection between the \mathcal{D} zeros of h_a and those of h_b .*

REMARK 3.3. Even if h_t crosses Σ some solutions may be able to be continued while others may not.

The (normalized) condition number [52] is a quantity describing “how close to singular” a zero is. Given $h \in \mathcal{H}_{(d)}$ and $z \in \mathbb{P}(\mathbb{C}^{n+1})$, let

$$(3.4) \quad \mu(f, z) = \|f\| \| (Dh(z) |_{z^\perp})^{-1} \text{Diag}(\|z\|^{d_i-1} d_i^{1/2}) \|_2,$$

and $\mu(f, z) = +\infty$ if $Dh(z) |_{z^\perp}$ is not invertible. Sometimes μ is denoted μ_{norm} or μ_{proj} but we prefer to keep the more simple notation here. One of the most important properties of μ is that it is an upper bound for the norm of the (locally defined) implicit function related to π_1 in (3.2). Namely, let $(\dot{h}, \dot{\zeta}) \in T_{(h, \zeta)} \mathcal{V}$ where $(h, \zeta) \in \mathcal{V}$ is such that $\mu(h, \zeta) < +\infty$. Then,

$$(3.5) \quad \|\dot{\zeta}\| \leq \mu(h, \zeta) \|\dot{h}\|, \quad \mu(h, \zeta) \geq \sqrt{n}.$$

We also have the following result.

THEOREM 3.4 (Condition Number Theorem, [52]).

$$\mu(h, \zeta) = \frac{1}{\sin(d_R(h, \Sigma_\zeta))},$$

where d_R is the Riemannian distance in $\mathbb{P}(\mathcal{H}_{(d)})$ and

$$\Sigma_\zeta = \{h \in \mathbb{P}(\mathcal{H}_{(d)}) : h(\zeta) = 0, \text{ and } Dh(\zeta) |_{\zeta^\perp} \text{ is not invertible}\}.$$

Note that this is a version of the classical Condition Number Theorem of linear algebra (see Theorem 6.5 below). The existence of approximate zeros in the sense of Definition 2.4 above is also guaranteed by this condition number, as was noted in [52]. More precisely:

THEOREM 3.5 (μ -Theorem, [52]). *There exists a constant $u_0 > 0$ ($u_0 = 0.17586$ suffices) with the following property. Let $(h, \zeta) \in \mathcal{V}$ and let $z \in \mathbb{P}(\mathbb{C}^{n+1})$ satisfy*

$$d_R(z, \zeta) \leq \frac{u_0}{d^{3/2} \mu(h, \zeta)}.$$

Then, z is an approximate zero of h with associated zero ζ .

4. The complexity of following a homotopy path

The sentence “can be continued” in the discussion of Section 3 can be made much more precise, by defining an actual path-following method. It turns out that the unique method that has actually been proved to correctly follow the homotopy paths and at the same time achieve some known complexity bound is the most simple one, which only uses the projective Newton operator, and not an ODE solver step.

PROBLEM 4.1. *It would be an interesting project to compare the overall cost of using a higher order ODE solver to the projective Newton-based method we describe below. Higher order methods or even predictor-corrector methods may require fewer steps but be more expensive at each step so a total cost comparison is in order. Some experience indicates that higher order methods are rarely cheaper, if ever. See [39, 40].*

More precisely, the projective Newton-based homotopy method is as follows. Given a C^1 path $\{h_t : a \leq t \leq b\} \subseteq \mathbb{P}(\mathcal{H}_{(d)})$, and given z_a an approximate zero of h_a with associated (exact) zero ζ_a , let $t_0 > 0$ be “small enough” and let

$$z_{a+t_0} = z_a - (Dh_{a+t_0}(z_a) |_{z_a^\perp})^{-1} h_{a+t_0}(z_a),$$

that is z_{a+t_0} is the result of one application of the projective Newton operator based on h_{a+t_0} to the point z_a . If z_a is an approximate zero of h_a and t_0 is small enough, then z_a can be close enough to the actual zero ζ_{a+t_0} of h_{a+t_0} to satisfy Theorem 3.5 and thus be an approximate zero of h_{a+t_0} as well. Then, by definition of approximate zero, z_{a+t_0} will be half-closer to ζ_{a+t_0} than z_a . This leads to an inductive process (choosing t_1 , then t_2 , etc. until h_b is reached) that, analysed in detail, can be made to work and actually programmed. The details on how to choose t_0 would take us too far apart from the topic, so we just give an intuitive explanation: if we are to move from (h_a, ζ_a) to $(h_{a+t_0}, \zeta_{a+t_0})$ we must be sure that we are far enough from Σ' to have our algorithm behaving properly. As the condition number essentially measures the distance to Σ' , it should be clear that the bigger the condition number, the smaller step t_0 we can take. This idea lead to the following result (see [56] for a weaker, earlier result):

THEOREM 4.2 ([51]). *Let $(h_t, \zeta_t) \subseteq \mathcal{V} \setminus \Sigma'$, $t \in [a, b]$ be a C^1 path. If the steps t_0, t_1, \dots are correctly chosen, then an approximate zero of h_b is reached at some point, namely there is a $k \geq 1$ such that $\sum_{i=0}^k t_i = b - a$ (k is the number of steps in the inductive process above.) Moreover, one can bound*

$$k \leq \lceil Cd^{3/2} L_\kappa \rceil,$$

where d is the maximum of the degrees in (d) , C is some universal constant, and

$$(4.1) \quad L_\kappa = \int_a^b \mu(h_t, \zeta_t) \|(\dot{h}_t, \dot{\zeta}_t)\| dt$$

is called the condition length of the path (h_t, ζ_t) . Moreover, the amount of arithmetic operations needed in each step is polynomial in the input size N , and hence the total complexity of the path-following procedure is a quantity polynomial in N and linear in L_κ .

There exist several ways to algorithmically produce the steps t_0, t_1, \dots in this theorem (and indeed the process has been programmed in two versions [12, 13]), but the details are too technical for this report, see [8, 27, 31]. We also point out that, if the path we are following is linear, i.e. $h_t = (1-t)h_0 + th_1$, and if the input coordinates are (complex) rational numbers, then all the operations can be carried out over the rationals without a dramatic increase of the bit size of intermediate results, see [13].

Note that since L_κ is a length it is independent of the C^1 parametrization of the path. If we specify a path of polynomial systems in $\mathcal{H}_{(d)}$ then we project the path of polynomials and solutions into \mathcal{V} to calculate the length. We may project from $\mathcal{H}_{(d)}$ to \mathbb{S} first and reparametrize if we wish. For example, we project the straight line segment $h_t = (1-t)g + th$ for $0 \leq t \leq 1$ into \mathbb{S} and reparametrize by arc-length. If $\|g\| = \|h\| = 1$ the resulting curve is

$$h_t = g \cos(t) + \frac{h - \langle h, g \rangle g}{\|h - \langle h, g \rangle g\|} \sin(t)$$

which is an arc of great circle through g and h . If $0 \leq t \leq d_R(g, h)$, then the arc goes from g to h . Here $d_R(g, h)$ is the Riemannian distance in \mathbb{S} between g and h which is the angle between them.

5. The problem of good starting points

We now come back to the original question in Smale's 17-th problem. Our plan is to analyse the complexity of an algorithm that we could call "linear homotopy": choose some $g \in \mathbb{S}$, $\zeta \in \mathbb{P}(\mathbb{C}^{n+1})$ such that $g(\zeta) = 0$ (we will call (g, ζ) a "starting pair"). For input $h \in \mathbb{S}$, consider the path contained in the great circle :

$$(5.1) \quad h_t = g \cos(t) + \frac{h - \langle h, g \rangle g}{\|h - \langle h, g \rangle g\|} \sin(t), \quad t \in [0, d_R(g, h)].$$

Then, use the method described in Theorem 4.2 to track how ζ_0 moves to $\zeta_{d_R(g, h)}$, a zero of $h_{d_R(g, h)} = h$, thus producing an approximate zero of h . We call this linear homotopy (maybe a more appropriate name would be "great circle homotopy") because great circles are projections on \mathbb{S} of segments in $\mathcal{H}_{(d)}$.

Assuming that the input h is uniformly distributed on \mathbb{S} , we can give an upper bound for the average number of arithmetic operations needed for this task (that is, the average complexity of the linear homotopy method) by a polynomial in N multiplied by the following quantity:

$$\frac{1}{\nu(\mathbb{S})} \int_{h \in \mathbb{S}} \int_0^{d_R(g, h)} \mu(h_t, \zeta_t) \|(\dot{h}_t, \dot{\zeta}_t)\| dt d\mathbb{S},$$

where h_t is defined by (5.1) and ζ_t is defined by continuation (the fact that $h_t \cap \Sigma = \emptyset$, and thus the existence of such ζ_t , is granted by Theorem 3.1 for most choices of (g, h)). It is convenient to replace this last expected value by a similar upper bound:

$$\mathcal{A}_1(g, \zeta) = \frac{1}{\nu(\mathbb{S})} \int_{h \in \mathbb{S}} \int_0^\pi \mu(h_t, \zeta_t) \|(\dot{h}_t, \dot{\zeta}_t)\| dt d\mathbb{S}.$$

Note that we are just replacing the integral from 0 to $d_R(g, h)$ by the distance from 0 to π .

We thus have:

THEOREM 5.1. *Let $(g, \zeta) \in \mathcal{V}$. The average complexity of linear homotopy with starting pair (g, ζ) is bounded above by a polynomial in N multiplied by $\mathcal{A}_1(g, \zeta)$.*

This justifies the following definition:

DEFINITION 5.2. Fix some polynomial¹ $p \in \mathbb{R}[x, y, z]$. We say that (g, ζ) is a good starting pair w.r.t. $p(x, y, z)$ if $\mathcal{A}_1(g, \zeta) \leq p(n, d, N)$ (which implies that the average number of steps of the linear homotopy is $O(d^{3/2}p(n, d, N))$.) From now on, if nothing is said, we assume $p(x, y, z) = \sqrt{2}\pi xz$. Thus, $(g, \zeta) \in \mathcal{V}$ is a good initial pair if $\mathcal{A}_1(g, \zeta) \leq \sqrt{2}\pi nN$.

So, if a good sequence of initial pair is known for all choices of n and the list of degrees (d) , then the total average complexity of linear homotopy is polynomial in N . In other words, finding good starting pairs for every choice of n and (d) gives a satisfactory solution to Problem (2.1).

In [56] the following pair² was conjectured to be a good starting pair (for some polynomial $p(x, y, z)$):

$$(5.2) \quad g(z) = \begin{cases} d_1^{1/2} z_0^{d_1-1} z_1, \\ \vdots \\ d_n^{1/2} z_0^{d_n-1} z_n \end{cases}, \quad \zeta = (1, 0, \dots, 0).$$

To this date, proving this conjecture is still an open problem. Some experimental data supporting this conjecture was shown in [12].

5.1. Choosing initial pairs at random: an Average Las Vegas algorithm for problem (2.1). One can study the average value of the quantity $\mathcal{A}_1(g, \zeta)$ described above. Most of the results in this section are based on the fact that the expected value of the square of the condition number is relatively small. This was first noted in [53], then this expected value was computed exactly in [16]:

THEOREM 5.3. *Let $h \in \mathbb{S}$ be chosen at random, and let ζ be chosen at random, with the uniform distribution, among the zeros of h . Then, the expected value of $\mu^2(h, \zeta)$ is at most nN . More exactly:*

$$\mathbb{E}_{h \in \mathbb{S}} \left(\frac{1}{\mathcal{D}} \sum_{\zeta: h(\zeta)=0} \mu(h, \zeta)^2 \right) = N \left(n \left(1 + \frac{1}{n} \right)^{n+1} - 2n - 1 \right) \leq nN.$$

¹Because $n, d \leq N$, we could just talk about a one variable polynomial $p(x)$ and change $p(n, d, N)$ to $p(N)$ in the following definition. However, we prefer here to be a bit more precise.

²The pair conjectured in [56] does not contain the extra $d_i^{1/2}$ factors. There is, however, some consensus that these extra factors should be added, for with these factors the condition number $\mu(g, \zeta) = n^{1/2}$ is minimal.

In particular, in the case of one homogeneous polynomial of degree d (i.e. $n = 1$,) we have:

$$\mathbb{E}_{h \in \mathbb{S}} \left(\sum_{\zeta: h(\zeta)=0} \mu(h, \zeta)^2 \right) = d(d+1).$$

Now we use some arguments which are very much inspired by ideas from integral geometry, one of the main contributions of Lluís Santaló to XX century mathematics. We can try to compute the expected value of $\mathcal{A}_1(g, \zeta)$. Although this can be done directly (see [18],) it is easier to first consider an upper bound of \mathcal{A}_1 : let us note from (3.5) that

$$(5.3) \quad \mathcal{A}_1(g, \zeta) \leq \frac{\sqrt{2}}{\nu(\mathbb{S})} \int_{h \in \mathbb{S}} \int_0^\pi \mu(h_t, \zeta_t)^2 dt d\mathbb{S}.$$

So, we have

$$\begin{aligned} \mathbb{E}_{g \in \mathbb{S}} \left(\sum_{\zeta: g(\zeta)=0} \mathcal{A}_1(g, \zeta) \right) &\leq \mathbb{E}_{g \in \mathbb{S}} \left(\sum_{\zeta: h(\zeta)=0} \frac{\sqrt{2}}{\nu(\mathbb{S})} \int_{h \in \mathbb{S}} \int_0^\pi \mu(h_t, \zeta_t)^2 dt d\mathbb{S} \right) = \\ &\sqrt{2} \mathbb{E}_{(g,h) \in \mathbb{S} \times \mathbb{S}} \left(\int_{f \in L_{g,h}} \sum_{\zeta: f(\zeta)=0} \mu(f, \zeta)^2 \right), \end{aligned}$$

where $L_{g,h}$ is the half-great circle in \mathbb{S} containing g, h , starting at g and going to $-g$ (we have to remove from this argument the case $h = -g$ but this is unimportant for integration purposes.) Note that we can define a measure and more generally a concept of integral in \mathbb{S} as follows: given any measurable function $q : \mathbb{S} \rightarrow [0, \infty)$, its integral is

$$(5.4) \quad \mathbb{E}_{(g,h) \in \mathbb{S} \times \mathbb{S}} \left(\int_{f \in L_{g,h}} q(f) \right).$$

Now, this last formula describes an invariant (with respect to the group of symmetries of \mathbb{S} , that can be identified with the unitary group of size $N + 1$ or with the orthogonal group of size $2N + 2$) measure in \mathbb{S} and is thus equal to a multiple of the usual measure in \mathbb{S} . In words, averaging over \mathbb{S} or over great circles in \mathbb{S} is the same, up to a constant. The constant is easy to compute by considering the constant function $q \equiv 1$. What we get is:

$$\mathbb{E}_{g \in \mathbb{S}} \left(\sum_{\zeta: g(\zeta)=0} \mathcal{A}_1(g, \zeta) \right) \leq \frac{\pi}{\sqrt{2}} \mathbb{E}_{h \in \mathbb{S}} \left(\sum_{\zeta: h(\zeta)=0} \mu(h, \zeta)^2 \right).$$

After this argument is made rigorous, we have (see [14, 15] for earlier versions of the following result:)

THEOREM 5.4 ([16]). *Let $g \in \mathbb{S}$ be chosen at random with the uniform distribution, and let ζ be chosen at random, with the uniform (discrete) distribution among the roots of g . Then, the expected value of $\mathcal{A}_1(g, \zeta)$ is at most $\frac{\pi}{\sqrt{2}} nN$. In particular, for such a randomly chosen pair (g, ζ) , with probability at least $1/2$ we have $\mathcal{A}_1(g, \zeta) \leq \sqrt{2} \pi nN$, that is, (g, ζ) is a good starting pair³.*

³Note that we are computing there the average of \mathcal{A}_1 not that of the integral of μ^2 as in [16]. From (5.3), the constant $\sqrt{2}$ has to be added to the formula in [16] in this context.

The previous result would be useless for describing an algorithm (because choosing a random zero of a randomly chosen $g \in \mathbb{S}$ might be a difficult problem) without the following one.

THEOREM 5.5 ([16]). *The process of choosing a random $g \in \mathbb{S}$ and a random zero ζ of g can be emulated by a simple linear algebra procedure.*

The details of the linear algebra procedure of Theorem 5.5 require the introduction of too much notation. We just describe the process in words: one has to choose a random $n \times (n + 1)$ matrix M with complex entries, compute its kernel (a projective point $\zeta \in \mathbb{P}(\mathbb{C}^{n+1})$) and consider the system $g \in \mathbb{S}$ that has ζ as a zero and whose linear part is given by M . A random higher-degree term has to be added to g , and then linear and higher-degree terms must be correctly weighted. This whole process has running time polynomial in N . We thus have:

COROLLARY 5.6. *The linear homotopy algorithm with the starting pair obtained as in Theorem 5.5 has average complexity⁴ $\tilde{O}(N^2)$.*

The word ‘‘average’’ in Corollary 5.6 must be understood as follows. For an input system h , let $T(h)$ be the expected running time of the linear homotopy algorithm, when (g, ζ) is randomly chosen following the procedure of Theorem 5.5. Then, the average value of $T(h)$ for random h is $\tilde{O}(N^2)$. This kind of algorithm is called Average Las Vegas, the ‘‘Las Vegas’’ term coming from the fact that a random choice has to be done. The user of the algorithm plays the role of a Las Vegas casino, not of a Las Vegas gambler: the chances of winning (i.e. getting a fast answer to our problem) are much higher than those of loosing (i.e. waiting for a long time before getting an answer.)

Some of the higher moments of $\mathcal{A}_1(g, \zeta)$ have also been proved to be small. For example, the second moment (thus, also the variance) of $\mathcal{A}_1(g, \zeta)$ is polynomial in N , as the following result shows:

THEOREM 5.7 ([18]). *Let $2 \leq k < 3$. Then, the expectation of $\mathcal{A}_1(g, \zeta)^k$ satisfies*

$$\mathbb{E}(\mathcal{A}_1(g, \zeta)^k) < \infty.$$

Moreover, let $2 \leq k < 3 - \frac{1}{2 \ln \mathcal{D}}$. Then, the expectation $\mathbb{E}(\mathcal{A}_1(g, \zeta)^k)$ satisfies,

$$\mathbb{E}(\mathcal{A}_1(g, \zeta)^k) \leq 2^{2k+k/2+4} e \pi^k n^{3k-4} N^2 \mathcal{D}^{4k-8} \ln \mathcal{D}.$$

In particular, $\mathbb{E}(\mathcal{A}_1(g, \zeta)^2) \leq 512e\pi^2 n^2 N^2 \ln \mathcal{D}$.

We have been concentrating on finding one zero of a polynomial system. But we could find k zeros $0 \leq k \leq \mathcal{D}$ by choosing k different random initial pairs using Theorem 5.5. This process is known from [16] to output every zero of the goal system h with the same probability $1/\mathcal{D}$, if $h \notin \Sigma$. Another option is to choose some initial system g which has k known zeros, and simultaneously continuing the k homotopy paths with the algorithm of Theorem 4.2. In the case of finding all zeros the sum of the number of steps to follow each path, is by Theorem 4.2 and (3.5), bounded above by a constant times

$$d^{3/2} \int_0^{d_R(g,h)} \sum_{\zeta_t: h_t(\zeta_t)=0} \mu(h_t, \zeta_t)^2 dt.$$

⁴We use here the $\tilde{O}(X)$ notation: this is the same as $O(X \log(X)^c)$ for some constant c , that is logarithmic factors are cleaned up to make formulas look prettier.

So for the great circle homotopies we have been discussing an analogue of Theorem 5.4 holds:

THEOREM 5.8 ([16]). *Let $g \in \mathbb{S}$ be chosen at random with the uniform distribution. Then, the expected value of $\int_0^{d_R(g,h)} \sum_{\zeta_t: h_t(\zeta_t)=0} \mu(h_t, \zeta_t)^2 dt$ is at most $\frac{\pi}{\sqrt{2}} nND$. In particular, for such a randomly chosen g , with probability at least $1/2$ we have $\int_0^{d_R(g,h)} \sum_{\zeta_t: h_t(\zeta_t)=0} \mu(h_t, \zeta_t)^2 dt \leq \sqrt{2} \pi nND$, that is, the linear homotopy for finding all zeros starting at g takes at most a constant times $d^{3/2} nND$ steps to output all zeros of h , on the average.*

Note that in general, one cannot write down all the \mathcal{D} zeros of g to begin with, so Theorem 5.8 does not immediately yield a practical algorithm.

We point out that, even for the case $n = 1$, no explicit descriptions of pairs (g, ζ) satisfying $\mathcal{A}_1(g, \zeta) \leq d^{O(1)}$ are known. Of course, no explicit polynomial $g \in \mathbb{S}$ is known in that case satisfying the claim of Theorem 5.8. An attempt to determine such a polynomial has led to some progress in the understanding of elliptic Fekete points, see Section 8.

5.2. The roots of unity combined with a method of Renegar: a quasi-polynomial time deterministic algorithm for problem (2.1). One can also ask for an algorithm for Problem (2.1) which does not rely on random choices (a deterministic algorithm). The search of a deterministic algorithm with polynomial running time for Problem (2.1) is still open, but a quasi-polynomial algorithm is known since [27].

This algorithm is actually a combination of two: on one hand, we consider the initial pair

$$(5.5) \quad g = \begin{cases} \frac{1}{\sqrt{2n}}(x_0^{d_1} - x_1^{d_1}) \\ \vdots \\ \frac{1}{\sqrt{2n}}(x_0^{d_1} - x_n^{d_n}) \end{cases}, \quad \zeta = (1, \dots, 1)$$

Then, we have:

THEOREM 5.9 ([27]). *The projective Newton-based homotopy method with initial pair (5.5) has average running time polynomial in N and n^d (recall that d is the maximum of the degrees).*

Theorem 5.9 is a consequence of the following stronger result:

THEOREM 5.10 ([27]). *The projective Newton-based homotopy method with initial pair $(g, \zeta) \in \mathcal{V}$ has average running time polynomial in N and in $\max\{\mu(g, \eta) : g(\eta) = 0\}$.*

Theorem 5.9 follows from Theorem 5.10 and the fact that $\mu(g, \eta) \leq 2(n+1)^d$ for g given by (5.5) for every zero η of g .

For small (say, bounded) values of d , the quantity n^d is polynomial in n and thus polynomial in N , but for big values of d the quantity n^d is not bounded by a polynomial in N , and thus Theorem 5.9 does not claim the existence of a polynomial running time algorithm. However, it turns out that there is a previously known algorithm, based on the factorization of the u -resultant, that has exponential running time for small degrees, but polynomial running time for high degrees (this

may seem contradictory, but it is not: when the degrees are very high, the input size is big, and thus bounding the running time by a polynomial in the input size is sometimes possible in this case.) More precisely:

THEOREM 5.11 ([27, 48]). *There is an algorithm with average running time polynomial in N and \mathcal{D} that, on input $h \in \mathbb{P}(\mathcal{H}_{(d)}) \setminus \Sigma$, outputs an approximate zero associated to every single exact zero of h .*

Note that \mathcal{D} is usually exponential in n , but as suggested above, if the degrees are very high compared to n , then \mathcal{D} can be bounded above by a polynomial in the input size N and thus the algorithm of Theorem 5.11 becomes a polynomial running time algorithm.

An appropriate combination of theorems 5.9 and 5.11, using the homotopy method of Theorem 5.9 for moderately low degrees and the symbolic-numeric method of Theorem 5.11 for moderately high degrees turns out to be quasipolynomial for every choice of n and (d) . Indeed:

THEOREM 5.12 ([27]). *The average (for random $h \in \mathbb{S}$) running time of the following procedure is $O(N^{\log \log N})$: on every input $h \in \mathbb{P}(\mathcal{H}_{(d)}) \setminus \Sigma$, run simultaneously the algorithms of theorems 5.9 and 5.11, stopping the computation whenever one of the two algorithms gives an output.*

Note that the running time of this algorithm is thus quasi-polynomial in N . Moreover, the algorithm is deterministic because it does not involve random choices.

5.3. Homotopy paths based on the evaluation at one point. Another approach to construct homotopies was considered in [57] and generalized in [4]. Given $h \in \mathcal{H}_{(d)}$ and $\zeta \in \mathbb{P}(\mathbb{C}^{n+1})$, consider $g = h - \hat{h}_\zeta$, where $\hat{h}_\zeta \in \mathcal{H}_{(d)}$ is defined as

$$\hat{h}_\zeta(z) = \text{Diag} \left(\frac{\langle z, \zeta \rangle^{d_i}}{\langle \zeta, \zeta \rangle^{d_i}} \right) h(\zeta).$$

Then, $g(\zeta) = 0$. So, we consider the homotopy $h_t = (1-t)g + th = h - (1-t)\hat{h}_\zeta$. We continue the zero ζ from $h_0 = g$ to $h_1 = h$. For any fixed ζ , for example $\zeta = e_0 = (1, 0, \dots, 0)$, the homotopy may be continued for almost all $h \in \mathcal{H}_{(d)}$. Let

$$K(h, \zeta) = \text{number of steps sufficient to continue } \zeta \text{ to a zero of } h,$$

and

$$K(h) = \mathbb{E}_{\zeta \in \mathbb{P}(\mathbb{C}^{n+1})}(K(h, \zeta)).$$

Then,

THEOREM 5.13 ([4]).

$$\mathbb{E}_{h \in \mathcal{H}_{(d)}}(K(h)) \leq \frac{Cd^{3/2}\Gamma(n+1)2^{n-1}}{(2\pi)^N\pi^n} \int_{h \in \mathcal{H}_{(d)}} \left(\sum_{\eta: h(\eta)=0} \frac{\mu(h, \eta)^2}{\|h\|^2} \Theta(h, \eta) \right) e^{-\|h\|^2/2} dh,$$

where

$$\Theta(h, \eta) = \int_{\zeta \in B(h, \eta)} \frac{(\|h\|^2 - T^2)^{1/2}}{T^{2n-1}} \Gamma(T^2/2, n) e^{T^2/2} d\zeta,$$

$$T = \|\text{Diag}(\|\zeta\|^{-d_i})h(\zeta)\|,$$

and $\Gamma(\alpha, n) = \int_\alpha^{+\infty} t^{n-1} e^{-t} dt$ is the incomplete gamma function.

In Theorem 5.13, $B(h, \eta)$ is the basin of η , which we now define. Suppose η is a non-degenerate zero of $h \in \mathcal{H}_{(d)}$. We define the basin of η , $B(h, \eta)$, as those $\zeta \in \mathbb{P}(\mathbb{C}^{n+1})$ such that the zero ζ of $g = h - \hat{h}_\zeta$ continues to η for the homotopy $h_t = (1-t)g + th$. We observe that the basins are open sets.

Not much is known about $E(K)$. See [4] for precise questions and motivations. Here is one:

PROBLEM 5.14. *Is $E(K)$ a quantity polynomial in N ?*

6. The condition Lipschitz–Riemannian structure

Let us now turn our sight back to (4.1). If we drop the condition number $\mu(h_t, \zeta_t)$ from that formula, we get

$$L = \int_a^b \|\dot{h}_t, \dot{\zeta}_t\| dt,$$

that is simply the length of the path (h_t, ζ_t) in the solution variety \mathcal{V} , taking on \mathcal{V} the natural metric: the one inherited from that of the product $\mathbb{P}(\mathcal{H}_{(d)}) \times \mathbb{P}(\mathbb{C}^{n+1})$. The formula in (4.1) can now be seen under a geometrical perspective: L_κ is just the length of the path (h_t, ζ_t) when \mathcal{V} is endowed with the conformal metric obtained by multiplying the natural one by the square of the condition number. Note that this new metric is only defined on $\mathcal{W} = \mathcal{V} \setminus \Sigma'$. We call this new metric the condition metric in \mathcal{W} . This justifies the name condition length we have given to L_κ . Theorem 4.2 now reads simply as follows: the complexity of following a homotopy path (h_t, ζ_t) is at most a small constant $cd^{3/2}$ times the length of (h_t, ζ_t) in the condition metric. This makes the condition metric an interesting object of study: which are the theoretical properties of that metric? given $p, q \in \mathcal{W}$, what is the condition length of the shortest path joining p and q ?

The first thing to point out is that μ is not a C^1 function, as it involves a matrix operator norm. However, μ is locally Lipschitz. Thus, the condition metric is not a Riemannian metric (usually, one demands smoothness or at least C^1 for Riemannian metrics,) but rather we may call it a Lipschitz–Riemannian structure. This departs from the topic of most available books and papers dealing with geometry of manifolds, but there are still some things we can say. It is convenient to take a tour to a slightly more general kind of problems; that's the reason for the following section.

6.1. Conformal Lipschitz–Riemann structures and self-convexity. Let M be a finite-dimensional Riemannian manifold, that is a smooth manifold with a smoothly varying inner product defined at the tangent space to each point $x \in M$, let us denote it $\langle \cdot, \cdot \rangle_x$. Let $\alpha : M \rightarrow [0, \infty)$ be⁵ a Lipschitz function, that is, there exists some constant $K \geq 0$ such that

$$|\alpha(x) - \alpha(y)| \leq K d_R(x, y), \quad \forall x, y \in M,$$

where $d_R(x, y)$ is the Riemannian distance from x to y . Then, consider on each point $x \in M$ the inner product $\langle \cdot, \cdot \rangle_{\alpha, x} = \alpha(x) \langle \cdot, \cdot \rangle_x$. Note that this need no longer be smoothly varying with x , for $\alpha(x)$ is just Lipschitz. We call such a structure

⁵The reader may have in mind the case $\alpha(h, \zeta) = \mu(h, \zeta)^2$ defined in $M = \mathcal{V}$.

a (conformal) Lipschitz–Riemannian structure in \mathcal{M} , and call it the α –structure. The condition length of a C^1 path $\gamma(t) \subseteq \mathcal{M}$, $a \leq t \leq b$, is just

$$L_\alpha(\gamma) = \int_a^b \|\dot{\gamma}(t)\|_{\alpha, \gamma(t)} dt = \int_a^b \langle \dot{\gamma}(t), \dot{\gamma}(t) \rangle_{\alpha, \gamma(t)}^{1/2} dt$$

The distance between any two points $p, q \in \mathcal{M}$ in this α –structure is defined as

$$(6.1) \quad d_\alpha(p, q) = \inf_{\gamma(t) \subseteq \mathcal{V}} L_\alpha(\gamma), \quad p, q \in \mathcal{M},$$

where the infimum is over all C^1 paths with $\gamma(0) = p, \gamma(1) = q$.

A path $\gamma(t)$, $a \leq t \leq b$ is called a minimizing geodesic if $L_\alpha(\gamma) = d_\alpha(\gamma(a), \gamma(b))$ and $\|\dot{\gamma}(t)\|_{\alpha, \gamma(t)} \equiv 1$, that is, if it minimizes the length of curves joining its extremal points and if it is parametrized by arc–length. Then, a curve $\gamma(t) \subseteq \mathcal{M}$, for t in some (possibly unbounded) interval I is called a geodesic if it is locally minimizing, namely if for every t in the interior of I there is some interval $[a, b] \subseteq I$ containing t and such that $\gamma|_{[a, b]}$ is a minimizing geodesic.

Each connected component of the set \mathcal{M} with the metric given by d_α is a path metric space, and it is locally compact because \mathcal{M} is a smooth finite–dimensional manifold. We are in a position to use Gromov’s version of the classical Hopf–Rinow theorem [36, Th.1.10], and we have:

THEOREM 6.1. *Let \mathcal{M} and α be as in the discussion above. Assume additionally that \mathcal{M} is connected and that (\mathcal{M}, d_α) is a complete metric space. Then:*

- each closed, bounded subset is compact,
- each pair of points can be joined by a minimizing geodesic.

Theorem 6.1 gives us sufficient conditions for conformal Lipschitz–Riemannian structures to be “well defined” in the sense that the infimum of (6.1) becomes a minimum. We can go further:

THEOREM 6.2 ([11]). *In the notation above, any geodesic is of class C^{1+Lip} , that is it is C^1 and has a Lipschitz derivative.*

See [22] for an early version of Theorem 6.2 and for experiments related to this problem.

One often thinks of the function α as some kind of “squared inverse of the distance to a bad set”, so for each connected component of \mathcal{M} the set (\mathcal{M}, d_α) will actually be complete.

A natural property to ask about is the following: given $p, q \in \mathcal{M}$, and given a geodesic $\gamma(t)$ such that $\gamma(a) = p, \gamma(b) = q$, does α attain its maximum on γ in the extremes? Namely, if we think on α as some kind of squared inverse to a bad set, do we have to get closer to the bad set than what we are in the extremes?

EXAMPLE 6.3. A model to think of is Poincaré half–plane with the metric given by the usual scalar product in $\mathbb{R}^2 \cap \{y > 0\}$, multiplied by $\alpha(x, y) = y^{-2}$. Geodesics then become just portions of vertical lines or half–circles with center at the axis $y = 0$. It is clear that, to join any two points, the geodesic does not need to become closer to the bad set $\{y = 0\}$.

We can ask for more: we say that α is self–convex (an abbreviation for self–log–convex) if for any geodesic $\gamma(t)$, the following is a convex function:

$$t \mapsto \log(\alpha(\gamma(t))).$$

Note that this condition is stronger than just asking for $t \mapsto \alpha(\gamma(t))$ to be convex, and thus stronger than asking for the maximum of α on γ to be at the extremal points.

6.2. Convexity properties of the condition number. We have the following result:

THEOREM 6.4 ([10]). *Let $k \geq 1$ and let $N \subseteq \mathbb{R}^k$ be a C^2 submanifold without boundary of \mathbb{R}^2 . Let $U \subseteq \mathbb{R}^n \setminus N$ be the biggest open set all of whose points have a unique closest point in N . Then, the function $\alpha(x) = \text{distance}(x, N)^{-2}$ is self-convex in U .*

Note that Theorem 6.4 is a more general version of Example 6.3, where the horizontal line $\{y = 0\}$ is changed to a submanifold N .

A well-known result usually attributed to Eckart and Young [35] and to Schmidt and Mirsky (see [61]) relates the usual condition number of a full rank rectangular matrix to the inverse distance to the set of rank-deficient matrices:

THEOREM 6.5 (Condition Number Theorem of linear algebra). *Let $A \in \mathbb{C}^{mn}$ be a $m \times n$ matrix for some $1 \leq m \leq n$. Let $\sigma_1(A), \dots, \sigma_m(A)$ be its singular values. Then,*

$$\sigma_m(A) = \text{distance}(A, \{\text{rank-deficient matrices}\}).$$

In particular, in the case of square maximal rank matrices, we can rewrite this as $\|A^{-1}\| = \text{distance}(A, \{\text{rank-deficient matrices}\})^{-1}$, that is the (unscaled) condition number $\|A^{-1}\|$ equals the inverse of the distance from A to the set of singular matrices. We more generally call $\sigma_m^{-1}(A)$ the unscaled condition number of a (possibly rectangular) full-rank matrix A .

One feels tempted to conclude from theorems 6.4 and 6.5 that the function sending a full-rank complex matrix A to the squared inverse of its smallest singular value (i.e. to the square of its unscaled condition number) should be self-convex. Indeed, one cannot apply Theorem 6.4 because the set of rank-deficient matrices is not a C^2 manifold, and because the distance to it is for many matrices (more precisely: whenever the multiplicity of the smallest singular value is greater than 1) not attained in a single point. It takes a considerable effort to prove that the result is still true:

THEOREM 6.6 ([11]). *The function defined in the space of full-rank $m \times n$ matrices, $1 \leq m \leq n$, as the squared inverse of the unscaled condition number, is self-convex.*

Note that this implies that, given any two complex matrices A, B of size $m \times n$, and given any geodesic $\gamma(t)$, $a \leq t \leq b$ in the α -structure defined in

$$\mathbb{C}^{mn} \setminus \{\text{rank-deficient matrices}\}$$

by $\alpha(C) = \sigma_m(C)^{-2}$ such that $\gamma(a) = A$, $\gamma(b) = B$, the maximum of α along γ is $\alpha(A)$ or $\alpha(B)$.

Note that, if a similar result could be stated for the α -structure defined by $(h, \zeta) \mapsto \mu(h, \zeta)^2$ in \mathcal{W} , we would have quite a nice description of how geodesics in the condition metric of \mathcal{W} are. Proving this is still an open problem:

PROBLEM 6.7. *Prove or disprove μ^2 is a self-convex function in \mathcal{W} .*

Note that from Theorem 3.4, the function μ^2 is not exactly the squared inverse of the distance to a submanifold, but it is still something similar to that. This makes it plausible to believe that Problem 6.7 has an affirmative answer. A partial answer is known:

THEOREM 6.8 ([11]). *The function $h \mapsto \mu^2(h, e_0)$ defined in the set $\{h \in \mathbb{P}(\mathcal{H}_{(d)}) : h(e_0) = 0\}$ is self-convex. Here, $e_0 = (1, 0, \dots, 0)$.*

7. Condition geodesics and the geometry of \mathcal{W}

Although we do not have an answer to Problem 6.7, we can actually state some bounds that give clues on the properties of the geodesics in the condition structure in \mathcal{W} . More precisely:

THEOREM 7.1 ([17]). *For every two pairs $(h_1, \zeta_1), (h_2, \zeta_2) \in \mathcal{W}$, there exists a curve $\gamma_t \subseteq \mathcal{W}$ joining (h_1, ζ_1) and (h_2, ζ_2) , and such that*

$$L_\kappa(\gamma_t) \leq 2cnd^{3/2} + 2\sqrt{n} \ln \left(\frac{\mu(h_1, \zeta_1)\mu(h_2, \zeta_2)}{n} \right),$$

c a universal constant.

In the light of Theorem 4.2, this means that if one can find geodesics in the condition structure in \mathcal{W} , one would be able to follow these paths in very few steps: just logarithmic in the condition number of the starting pair and the goal pair.

COROLLARY 7.2. *A sufficient number of projective Newton steps to follow some path in \mathcal{W} starting at the pair (g, e_0) of (5.2) to find an approximate zero associated to a solution ζ of a given system $h \in \mathbb{P}(\mathcal{H}_{(d)})$ is*

$$cd^{3/2} \left(nd^{3/2} + \sqrt{n} \ln \left(\frac{\mu(h, \zeta)}{\sqrt{n}} \right) \right),$$

c a universal constant.

Note that only the logarithm of the condition number appears in Corollary 7.2. Thus, if one could find an easy way to describe condition geodesics in \mathcal{W} , the average complexity of approximating them using Theorem 4.2 would involve just the expectation of the average of $\ln(\mu)$, not that of μ^2 as in Theorem 5.3. As a consequence, the average number of steps needed by such an algorithm would be $O(nd^3 \ln N)$. See [18, Cor. 3] for a more detailed statement of this fact. At this point we ask a rather naive, vague question:

PROBLEM 7.3. *May homotopy methods be useful in solving linear systems of equations? Might using geodesics help as in Corollary 7.2 and the comments above?*

Large sparse systems are frequently solved by iterative methods and the condition number plays a role in the error estimates. So Problem (7.3) has some plausibility.

REMARK 7.4. There is an exponential gap between the average number of steps needed by linear homotopy $O(d^{3/2}nN)$ and those promised by the condition geodesic-based homotopy (which stays at a theoretical level by now, because one cannot easily describe those geodesics). This exponential gap occurs frequently in theoretical computer science. For example NP-complete problems are solvable in

simply exponential time but polynomial with a witness. The estimates for homotopies with condition geodesics may likely serve as a lower bound for what can be achieved. Also, properties of geodesics as we learn them can inform the design of homotopy algorithms.

There is more we can say about the geometry (and topology) of \mathcal{W} , by studying the Frobenius condition number in W , which is defined as follows:

$$\tilde{\mu}(h, \zeta) = \|h\| \|Dh(\zeta)^\dagger \text{Diag}(\|\zeta\|^{d_i-1} d_i^{1/2})\|_F, \quad \forall (h, \zeta) \in W,$$

where $\|\cdot\|_F$ is Frobenius norm (i.e. $\text{Trace}(L^*L)^{1/2}$ where L^* is the conjugate transpose of L) and † is Moore-Penrose pseudoinverse.

REMARK 7.5. The Moore-Penrose pseudoinverse $L^\dagger : \mathbb{F} \rightarrow \mathbb{E}$ of a linear operator $L : \mathbb{E} \rightarrow \mathbb{F}$ of finite dimensional Hilbert spaces is defined as the composition

$$(7.1) \quad L^\dagger = i_{\mathbb{E}} \circ (L|_{\text{Ker}(L)^\perp})^{-1} \circ \pi_{\text{Image}(L)},$$

where $\pi_{\text{Image}(L)}$ is the orthogonal projection on image L , $\text{Ker}(L)^\perp$ is the orthogonal complement of the nullspace of L , and $i_{\mathbb{E}}$ is the inclusion. If A is a $m \times (n+1)$ matrix and $A = UDV^*$ is a singular value decomposition of A , $D = \text{Diag}(\sigma_1, \dots, \sigma_k, 0, \dots, 0)$ then we can write

$$(7.2) \quad A^\dagger = VD^\dagger U^*, \quad D^\dagger = \text{Diag}(\sigma_1^{-1}, \dots, \sigma_k^{-1}, 0, \dots, 0).$$

In [19] we prove that $\tilde{\mu}$ is an equivariant Morse function defined in \mathcal{W} with a unique orbit of minima given by the orbit \mathcal{B} of the pair of (5.2) under the action of the unitary group $(U, (h, \zeta)) \mapsto (h \circ U^*, U\zeta)$.

The function $\mathcal{A}_1(g, \zeta)$ or even its upper bound (up to a $\sqrt{2}$ factor) estimate ⁶

$$\mathcal{B}_1(g, \zeta) = \frac{1}{\nu(\mathbb{S})} \int_{h \in \mathbb{S}} \int_0^\pi \mu(h_t, \zeta_t)^2 dt d\mathbb{S}$$

is an average of μ^2 in great circles. This remark motivates the following

PROBLEM 7.6. *Is $\mathcal{A}_1(g, \zeta)$ or $\mathcal{B}_1(g, \zeta)$ also an equivariant Morse function whose only critical point set is a unique orbit of minima?*

If so, due to symmetry considerations, it is the orbit through the conjectured good starting point (5.2). Here, one may want to replace the condition number μ in the definition of \mathcal{B}_1 with a smooth version such as the Frobenius condition number. A positive solution to this problem solves our main problem: the conjectured good initial pair (5.2) is not only good but even best.

Because the Frobenius condition number is an equivariant Morse function, the homotopy groups of \mathcal{W} are equal to those of \mathcal{B} , that can be studied with standard tools from algebraic topology. In the case that $n > 1$, for example, we get:

$$\begin{aligned} \pi_0(\mathcal{W}) &= \{0\} \\ \pi_1(\mathcal{W}) &= \mathbb{Z}/a\mathbb{Z} \\ \pi_2(\mathcal{W}) &= \mathbb{Z} \\ \pi_3(\mathcal{W}) &= \pi_k(\mathcal{SU}_{n+1}) \quad (k \geq 3), \end{aligned}$$

where \mathcal{SU}_{n+1} is the set of special unitary matrices of size $n+1$, $a = \text{gcd}(n, d_1 + \dots + d_n - 1)$ and $\mathbb{Z}/a\mathbb{Z}$ is the finite cyclic group of a elements.

⁶see (5.3).

In particular, we see that if all the d_i 's are equal then $a = 1$ and \mathcal{W} is simply connected; in particular, any curve can be continuously deformed into a minimizing geodesic. See [19] for more results concerning the geometry of \mathcal{W} . We can also prove a lower bound similar to the upper bound of Theorem 7.1:

THEOREM 7.7. *let $\alpha : [a, b] \rightarrow \mathcal{W}$ be a C^1 curve. Then, its condition length is at least*

$$\frac{1}{d^{3/2}\sqrt{n+1}} \left| \ln \left(\frac{\mu(\alpha(a))}{\mu(\alpha(b))} \right) - \ln \sqrt{n+1} \right|.$$

REMARK 7.8. We have written Theorem 7.7 using the condition metric as defined in this paper. The original result [19, Prop. 11] was written for the so-called smooth condition length, obtained by changing μ to $\tilde{\mu}$ in the definition of the condition length. This change produces the $\sqrt{n+1}$ factors in Theorem 7.7.

In his article [59], Smale suggests that the input size of an instance of a numerical analysis problem should be augmented by $\log W(y)$ where $W(y)$ is a weight function "... to be chosen with much thought..." and he suggests that "the weight is to resemble the reciprocal of the distance to the set of ill-posed problems." That is the case here. The condition numbers we have been using are comparable to the distance to the ill-posed problems and figure in the cost estimates. It would be good to develop a theory of computation which incorporates the distance to ill-posedness, or condition number and distance to ill-posedness in case they may not be comparable, (and precision in the case of round-off error) more systematically so that a weight function will not require additional thought. For the case of linear programming Renegar [49] accomplished this. It is our main motivating example as well as the work we have described on polynomial systems. The book [28] is the current state of the art. The geometry of the condition metric will to our mind intervene in the analysis. If floating point arithmetic is the model of arithmetic used then ill-posedness will include points where the output is zero as well as points where the output is not Lipschitz.

8. The univariate case and elliptic Fekete points

Let us now center our attention in the univariate case, that, once homogenized, is the case of degree d homogeneous polynomials in two variables. Then,

$$\mu(h, \zeta) = d^{1/2} \|h\| \| (Dh(\zeta) |_{\zeta^\perp})^{-1} \| \zeta \|^{d-1}.$$

If we are given a univariate polynomial $f(x)$ and a complex zero z of f , we can also use the following more direct (and equivalent) formula for $\mu(h, \zeta)$ where h is the homogeneous counterpart of f and $\zeta = (1, z)$:

$$\mu(h, \zeta) = \frac{d^{1/2} (1 + |z|^2)^{\frac{d-2}{2}}}{|f'(z)|} \|h\|.$$

It was noted in [54] that the condition number is related to the classical problem of finding elliptic Fekete points, which we recall now in its computational form (see [9] for a survey on the state of art of this problem.)

Given d different points $x_1, \dots, x_d \in \mathbb{R}^3$, let $X = (x_1, \dots, x_d)$ and

$$\mathcal{E}(X) = \mathcal{E}(x_1, \dots, x_d) = - \sum_{i < j} \log \|x_i - x_j\|$$

be its logarithmic potential. Sometimes $\mathcal{E}(X)$ is denoted by $\mathcal{E}_0(X)$, $\mathcal{E}(0, X)$ or $V_N(X)$. Let $S(1/2)$ be the Riemann sphere in \mathbb{R}^3 , that is the sphere of radius $1/2$ centered at $(0, 0, 1/2)$, and let

$$m_d = \min_{x_1, \dots, x_d \in S(1/2)} \mathcal{E}(x_1, \dots, x_d)$$

be the minimum value of \mathcal{E} . A minimising d -tuple $X = (x_1, \dots, x_d)$ is called a set of elliptic Fekete points ⁷.

The computational problem of finding elliptic Fekete points is another of the problems in Smale's list ⁸.

Smale's 7th problem [60]: Can one find $X = (x_1, \dots, x_d)$ such that

$$(8.1) \quad \mathcal{E}(X) - m_d \leq c \log d, \quad c \text{ a universal constant.}$$

The first clue that this problem is hard comes from the fact that the value of m_d is not known, even to $O(d)$. A general technique (valid for Riemannian manifolds) given by Elkies shows that

$$m_d \geq \frac{d^2}{4} - \frac{d \log d}{4} + O(d).$$

Wagner [64] used the stereographic projection and Hadamard's inequality to get another lower bound. His method was refined by Rakhmanov, Saff and Zhou [45], who also proved an upper bound for m_d using partitions of the sphere. The lower bound was subsequently improved upon by Dubickas and Brauchart [34], [24]. The following result summarizes the best known bounds:

THEOREM 8.1. *Let C_d be defined ⁹ by*

$$m_d = \frac{d^2}{4} - \frac{d \log d}{4} + C_d d.$$

Then,

$$-0.4375 \leq \liminf_{d \rightarrow \infty} C_d \leq \limsup_{d \rightarrow \infty} C_d \leq -0.3700708\dots$$

The relation of this problem to the condition number relies on the fact that sets of elliptic Fekete points are naturally "well separated", and are thus good candidates to be the zeros of a "well-conditioned" polynomial, that is a polynomial all of whose zeros have a small condition number. In [54] Shub and Smale proved the following relation between the condition number and elliptic Fekete points.

THEOREM 8.2 ([54]). *Let $\zeta_1, \dots, \zeta_d \in \mathbb{P}(\mathbb{C}^2)$ be a set of projective points, and consider them as points in the Riemann sphere $S(1/2)$ with the usual identification $\mathbb{P}(\mathbb{C}^2) \cong S(1/2)$. Let h be a degree d homogeneous polynomial such that its zeros are ζ_1, \dots, ζ_d . Then,*

$$\max\{\mu(h, \zeta_i) : 1 \leq i \leq d\} \leq \sqrt{d(d+1)} e^{\mathcal{E}(\zeta_1, \dots, \zeta_d) - m_d}.$$

⁷Such a d -tuple can also be defined as a set of d points in the sphere which maximize the product of their mutual distances.

⁸Smale thinks on points in the unit sphere, but we may think on points in the Riemann sphere, as the two problems are equivalent by sending $(a, b, c) \in S(1/2)$ to $2(a, b, c) - (0, 0, 1)$.

⁹The result in the original sources is written for the unit sphere, we translate it here to the Riemann sphere.

In particular, is x_1, \dots, x_d are a set of elliptic Fekete points, then

$$\max\{\mu(h, \zeta_i) : 1 \leq i \leq d\} \leq \sqrt{d(d+1)}.$$

REMARK 8.3. Let \Re and \Im be, respectively, the real and complex part of a complex number. Here is alternative, equivalent definition for h and the ζ_i . Instead of considering projective points in $\mathbb{P}(\mathbb{C}^2)$ we may just consider a set of complex numbers $z_1, \dots, z_d \in \mathbb{C}$. Then, for $1 \leq i \leq d$, we can define $\zeta_i \in \mathbb{S}$ as

$$(8.2) \quad \zeta_i = \left(\frac{\Re(z_i)}{1 + |z_i|^2}, \frac{\Im(z_i)}{1 + |z_i|^2}, \frac{1}{1 + |z_i|^2} \right)^T \in S(1/2), \quad 1 \leq i \leq d,$$

f as the polynomial whose zeros are z_1, \dots, z_d , and h as the homogeneous counterpart of f .

There exists no explicit known way of describing a sequence of polynomials satisfying $\max\{\mu(h, \zeta) : h(\zeta) = 0\} \leq d^c$, for any fixed constant c and $d \geq 1$. Theorem 8.2 implies that, if a d -tuple satisfying (8.1) can be described for any d , then such a sequence of polynomials can also be generated. From Theorem 5.10, such h 's are good starting points for the linear homotopy method, both for finding one root and for finding all roots. So, solving the elliptic Fekete points problem solves the starting point problem for $n = 1$. The reciprocal question is: does solving the starting point problem for $n = 1$ help with the Fekete point problem?

PROBLEM 8.4. Suppose $n = 1$ and $g \in \mathbb{S}$ minimizes $\sum_{\zeta: g(\zeta)=0} \mu(g, \zeta)^2$. Do ζ_1, \dots, ζ_d (the zeros of g , seen as points in $S(1/2)$) solve Smale's 7-th problem?

We have seen in Theorem 5.3 that the condition number of (h, ζ) where h is chosen at random and ζ is uniformly chosen at random among the zeros of h , grows polynomially in d . Then, Theorem 8.2 suggests that spherical points associated with zeros of random polynomials might produce small values of \mathcal{E} . We can actually put some numbers to this idea. First, one can easily compute the average value of \mathcal{E} when x_1, \dots, x_d are chosen at random in $S(1/2)$, uniformly and independently with respect to the probability distribution induced by Lebesgue measure in $S(1/2)$:

$$E_{X \in S(1/2)^d} \mathcal{E}(X) = \frac{d^2}{4} - \frac{d}{4}.$$

By comparing this with Theorem 8.1, we can see that random choices of points in the sphere already produce pretty low values of the minimal energy. One can prove that random polynomials actually produce points which behave better with respect to \mathcal{E} :

THEOREM 8.5 ([3]). Let $n = 1$ and $h \in \mathbb{S}$ be chosen at random w.r.t. the uniform distribution in \mathbb{S} . Let $\zeta_1, \dots, \zeta_d \in S(1/2)$ be the zeros of h . Then, the expected value of $\mathcal{E}(\zeta_1, \dots, \zeta_d)$ equals

$$\frac{d^2}{4} - \frac{d \log d}{4} - \frac{d}{4}.$$

By comparing this with Theorem 8.1, we conclude that spherical points coming from zeros of random polynomials agree with the minimal value of \mathcal{E} , to order $O(d)$.

This result fits into a more general¹⁰ result related to random sections on Riemann surfaces, see [65, 66].

9. The algebraic eigenvalue problem

The double fibration scheme proposed in (3.2) has been – at least partly – successfully used in other contexts. For example, in [1] a similar projection scheme (9.1)

$$\begin{array}{ccc} & \mathcal{V}_{\text{eig}} & = \{((A, \lambda), v) \in \mathbb{P}(\mathbb{C}^{n^2+1}) \times \mathbb{P}(\mathbb{C}^n) : Av = \lambda v\} \\ \pi_1 \swarrow & & \searrow \pi_2 \\ \mathbb{P}(\mathbb{C}^{n^2+1}) & & \mathbb{P}(\mathbb{C}^{n+1}) \end{array}$$

was used to study the complexity of a homotopy-based eigenvalue algorithm, obtaining the following:

THEOREM 9.1. *A homotopy algorithm can be designed that continues an eigenvalue–eigenvector pair (λ_0, v_0) of a $n \times n$ matrix A_0 to one (λ_1, v_1) of another matrix A_1 , the number of steps bounded above by*

$$c \int_0^1 \|(\dot{A}, \dot{\lambda}, \dot{v})\| \mu_{\text{eig}}(A, \lambda, v) dt,$$

c a universal constant. Here, μ_{eig} is the condition number¹¹ for the algebraic eigenvalue problem, defined as

$$(9.2) \quad \mu_{\text{eig}}(A, \lambda, v) = \max \left\{ 1, \|A\|_F \|\pi_{v^\perp}(\lambda I_n - A)|_{v^\perp}^{-1}\| \right\},$$

*where $\|A\|_F = \text{trace}(A^*A)^{1/2}$ is the Frobenius norm of A .*

Of course, we do not intend to summarize here the enormous amount of methods and papers dealing with the eigenvalue problem (see [61] for example). We just point out that there exists no proven polynomial-time algorithm for approximating eigenvalues (although different numerical methods achieve spectacular results in practice.) See [44] for some statistics about the QR (and Toda) algorithms for symmetric matrices. We don't know a good reference for the more difficult general case. Unshifted QR is not the fast algorithm of choice. The QR algorithm with Francis double shift executed on upper Hermitian matrices should be the gold standard.

PROBLEM 9.2. *Does the QR algorithm with Francis double shift fail to attain convergence on an open subset of upper Hessenberg matrices?*

See [6] for open sets where Rayleigh quotient iteration fails, and [5] for a proof of convergence for normal matrices as well as a good introduction to the dynamics involved.

Theorem 9.1 can probably be used in an analysis similar to that of Section 5 to complete a complexity analysis. Note that the integral in Theorem 9.1 is very similar in spirit to that of (4.1). This allows to introduce a condition metric in

¹⁰Steve Zelditch tells us that “the relation between the special case of the round metric on $S(1/2)$ and the general metric on any Riemann surface is that the expansion terminates on $S(1/2)$ because the Fubini-Study metric is balanced, i.e. the szego kernel is constant on the diagonal. For general metrics it will not terminate.”

¹¹A quantity similar in spirit to the condition number μ for the polynomial system solving problem.

\mathcal{V}_{eig} . Some of the results in previous sections can be adapted to this new case. For example, an analogue of Theorem 7.1 holds (i.e. short geodesics exist,) see [2].

The eigenvalue problem and the problem of finding roots of a polynomial in one variable are, of course, connected. Given an $n \times n$ matrix A we may compute the characteristic polynomial of A , $p(z) = \det(zI - A)$ and then solve $p(z)$. The zeros of $p(z)$ are the eigenvalues of A . Trefethen and Bau [62] write “This algorithm is not only backward unstable but unstable and should not be used”. Indeed when presented with a univariate polynomial $p(z)$ to solve, numerical linear algebra packages may convert the problem to an eigenvalue problem by considering the companion matrix of $p(z)$ and then solve the eigenvalue problem. If $p(z) = z^d + a_{d-1}z^{d-1} + \dots + a_0$ the companion matrix is

$$\begin{pmatrix} 0 & 0 & 0 & \cdots & 0 & -a_0 \\ 1 & 0 & 0 & \cdots & 0 & -a_1 \\ 0 & 1 & 0 & \cdots & 0 & -a_2 \\ \vdots & & \ddots & \ddots & \vdots & \vdots \\ \vdots & & \ddots & \ddots & 0 & -a_{d-2} \\ 0 & \cdots & \cdots & 0 & 1 & -a_{d-1} \end{pmatrix},$$

which is already in upper Hessenberg form. So conceivably Francis double shifted QR may fail to converge on an open set of companion matrices?

Let us recall that the condition number of a polynomial and root is a property of the output map as a function of the input. So it doesn't depend on the algorithms to solve the problem. This motivates the following

PROBLEM 9.3. *What might explain the experience of numerical analysts, relating the polynomial solving methods versus that of eigenvalue solving? Might the condition number of the eigenvalue problem have small average over the set of $n \times n$ matrices with a given characteristic polynomial?*

Finally, we can consider the problem $Av = \lambda v$ as a system of n quadratic equations in n unknowns. By Bezout's theorem, after we homogenize, we expect 2^n roots counted with multiplicity. But there are only n eigenvalues. In [1, 2] it is shown that the use of multihomogeneous Bezout theorem yields the correct zero count for this problem. Thus, a reasonable thing to do is to introduce a new variable α and consider the bilinear equation $A\alpha v = \lambda v$ which is bilinear in (α, λ) and v .

PROBLEM 9.4 (see [32]). *Prove an analogue of Theorem 9.1 in the general multihomogeneous setting.*

Appendix A. A model of computation for machines with round-off and input errors

This section has been developed in discussions with Jean Pierre Dedieu and his colleagues Paola Boito and Guillaume Chèze. We thank Felipe Cucker for helpful comments.

A.1. Introduction. During the second half of the 20th century, with the emergence of computers, algorithms have taken a spectacular place in mathematics, especially numerical algorithms (linear algebra, ode's, pde's, optimization), but

also symbolic computation. In this context, complexity studies give a better understanding of the intrinsic difficulty of a problem, and describe the performance of algorithms which solve such problems. One can associate the classical Turing model to symbolic computation based on integer arithmetic, and the BSS model to scientific computation on real numbers. However this ideal picture suffers from an important defect. Scientific computation does not use the exact arithmetic of real numbers but floating-point numbers and a finite precision arithmetic. Thus, a numerical algorithm designed on real numbers and the same algorithm running in finite precision arithmetic give a priori two different results. Any numerical analysis undergraduate book has at least one chapter dealing with the precision of numerical computations. See for example [62] or [38]. Yet, there is no solid approach to the definition and study of a model of computation including this aspect, as well as the role that conditioning of problems should play in the complexity estimates.

Besides linear algebra problems and iterative processes, a key point to bear in mind is that we sometimes use floating point computers to answer decision (i.e Yes/No) problems, as is this matrix singular? or does this polynomial have a real zero?. The first attempts to use round-off machines to study decision problems are [30], and [29]. The authors consider questions like: under which conditions is the decision taken by a BSS machine the same as the decision taken by the corresponding round-off machine? Or, under which conditions is the decision taken by a round-off machine for a given input the same as the decision taken by the BSS machine on a nearby input?

In these pages we point towards the development of a theory of finite precision computation via a description of round-off machines, size of an input, cost of a computation, single (resp. multiple) precision computations (a computation is “single precision” when a sufficient round-off unit δ to reach relative precision u for any input x in the considered range is proportional to u), finite precision computability and finite precision decidability. These concepts have to be related to the intrinsic characteristics of the problem: its condition number (the local Lipschitz constant of the solution map), and its posedness (the distance to ill-posed problems).

The model we propose is inspired by the BSS model but it stays close to real-life numerical computation. We prefer relative errors to absolute ones (this is the basis of the usual floating-point arithmetic.) We mention two papers of interest about the foundations of scientific computing, [25, 26], with a point of view different than ours.

A.2. Round-off machine. A round-off machine is an implementation of a BSS-machine accounting for input error and round-off error of computations. These errors may mimic a particular floating point arithmetic but are designed to be more general. In particular, they are not tied down to a particular floating point model. Let \mathbb{R}^∞ be the disjoint union of the sets \mathbb{R}^n , $n \geq 0$. For given $x \in \mathbb{R}^\infty$ we define $\|x\| = \max_i |x_i|$. A subset $U \subseteq \mathbb{R}^\infty$ is open if it is the disjoint union of U_n with $U_n \subseteq \mathbb{R}^n$ an open set. For this topology, a mapping $f : \mathbb{R}^\infty \rightarrow \mathbb{R}$ is continuous iff each restriction $f_n = f|_{\mathbb{R}^n}$ is continuous.

A (real number) BSS machine M is a directed graph with with several kinds of nodes including an input node, with input $x \in \mathbb{R}^\infty$, output nodes, computation nodes where rational functions are generally computed but here we restrict ourselves without loss of generality to the standard arithmetic operations, branching nodes (we branch on an inequality of the form $y \geq 0$.) A machine is a decision machine

when the output is -1 or 1 . The halting set \mathcal{H} of M is the set of inputs giving rise to an output. We denote by $\mathcal{O} : \mathcal{H} \rightarrow \mathbb{R}^\infty$ the output map. There are a few technical concepts (mainly the input map $I_M(x)$ and the computing endomorphism H_M) associated to M , the unfamiliar reader may find formal definitions in [20, Chapters 2 and 3].

Given a BSS machine M with nodes $\{1, \dots, N\}$ and state space \mathbb{R}_∞ , we augment the state space \mathbb{R}_∞ by an extra copy of \mathbb{R} so the new state space is $\mathbb{R} \times \mathbb{R}_\infty$. The state space component of the input map is $(1, I_M(x))$. We define a new next node next state map \hat{H}_M by

$$\hat{H}_M(\eta, k, x) = (\pi_1(H_M(\eta, x)), k + 1, \pi_2(H_M(\eta, x))),$$

so the first coordinate acts as a counter (of the number of nodes visited by M). We say that the machine defined \hat{H}_M is a *counting BSS machine*. A little programming shows that adding this extra coordinate does nothing to change the computability or complexity theory of real BSS machines (indeed because $\mathbb{R} \times \mathbb{R}_\infty \equiv \mathbb{R}_\infty$, one can easily see that this newly defined machine is actually a BSS machine). We will moreover assume that our BSS machines are *elementary*, that is that the computation nodes of our machines contain only elementary operations, that is operations of the form $a \circ b$ where $a, b \in \mathbb{R}$ and $\circ \in \{+, -, \times, /\}$. It is a routine task to convert any given BSS machine into a counting elementary machine (this process can be done in many ways, because there are many different ways to compute a polynomial).

DEFINITION A.1 (Round-off machine associated to a given BSS machine). Given a counting, elementary BSS machine M defined over the real numbers and $0 \leq \delta \leq 1$, a round-off machine associated to M and δ is another machine (i.e. a directed graph with the same type of nodes as a BSS machine) denoted (M, δ) . The nodes and state space of (M, δ) are the same as for M . The input map $I_{(M, \delta)}$ of (M, δ) satisfies $|I_{(M, \delta)}(x)_j - I_M(x)_j| < \delta |I_M(x)_j|$ that is to say the relative error of the input is less than δ for every coordinate j . The next node next state of (M, δ) at a computation node has the same next node component as H_M , and the j th components of the next states satisfy $|H_{(M, \delta), state}(x)_j - H_{M, state}(x)_j| < \delta |H_{M, state}(x)_j|$, unless $H_{M, state}(x)_j = x_j$ in which case there is no error (i.e. $H_{(M, \delta), state}(x)_j = x_j$). The next node next state map is unchanged at a branch node or at a shift node.

Given any BSS machine M defined over the real numbers and $0 \leq \delta \leq 1$, a round-off machine associated to M and δ is a round-off machine (\tilde{M}, δ) associated to \tilde{M} and δ where \tilde{M} is some counting, elementary version of M .

REMARK A.2. The rounding error introduced at each computation node depends on the whole state and, because M is assumed to be a counting machine, the rounding error may thus depend on the counter. Thus, the rounding error introduced at a given node visited twice may be different (because the counter may be different). Note that the counter is also affected by rounding errors.

Note that a round off machine is not necessarily a BSS machine, and that given M and δ , there are many machines satisfying this definition. For example, M itself satisfies this definition for every δ . The power of the definition is that certain claims will hold for every such a round-off machine, allowing us to use just the defining properties and not the particular structure of a given round-off machine.

DEFINITION A.3. Given a BSS machine M and $0 < \delta < 1$, a δ pseudo-computation with input x is the sequence of pairs $(node, state)$ generated by *some* round-off machine associated to some counting, elementary version of M .

We also point out that not every BSS machine can be (reasonably) converted into a round-off machine. For example, assume that a BSS machine performs the operation $x = (x_1, \dots, x_N) \mapsto x_1 + x_N$. This machine must contain a loop counting up to N . If the form of the *if* node defining the loop is $k \geq 0$ (k the counter which is, say, diminished by 1 at each step) then an arbitrarily small error in the counter of the loop may produce that an associated round-off machine on input x outputs $x_1 + x_{N-1}$ instead of $x_1 + x_N$. A clear way out is to consider the slightly different BSS machine which checks if $-1/2 \leq k \leq 1/2$ instead of $k \geq 0$. Then, a round-off machine with reasonable precision $\delta = O(1/N)$ will do the job. Note that this fits perfectly into the definition of single precision computation (A.7) below. This also reflects the fact, known to every numerical analyst or programmer, that not every program is suitable for floating point conversion: a little care needs to be taken!

A.3. Computability. In the sequel, we will only consider functions $f : \Omega \subseteq \mathbb{R}^\infty \rightarrow \mathbb{R}^\infty$ such that, for each n , the restriction f_n of f to $\Omega_n = \Omega \cap \mathbb{R}^n$ takes its values in \mathbb{R}^m for an m depending only on n .

Such a function is round-off computable when there exists a BSS machine M such that for any $x \in \Omega$ and any $0 < \epsilon < 1$, there exists a $\delta(x, \epsilon)$ such that any round-off machine $(M, \delta(x, \epsilon))$ outputs $\tilde{O}(x)$ with

$$|\tilde{O}(x)_j - f(x)_j| \leq \epsilon |f(x)_j|,$$

that is the output of $(M, \delta(x, \epsilon))$ is coordinatewise equal to $f(x)$ up to relative error ϵ . Equivalently, we say that M round-off computes f if given $x \in \Omega$ and $0 < \epsilon < 1$, there is $\delta(x, \epsilon)$ such that all $\delta(x, \epsilon)$ pseudo-computations of M on input x output $f(x)$ with relative precision ϵ .

EXAMPLE A.4. The function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$, $f(x, y) = xy$ (we can let it be zero in $\mathbb{R}^\infty \setminus \mathbb{R}^2$) is round-off computable. Indeed, let $x, y \neq 0$ and $0 < \epsilon < 1$. The output of a round-off machine (M, δ) associated to the natural BSS machine for computing $f(x, y)$ is a number

$$z = xy(1 + \delta_1)(1 + \delta_2)(1 + \delta_3),$$

for some $\delta_1, \delta_2, \delta_3$ bounded in absolute value by δ . It is useful to note the elementary inequality

$$(A.1) \quad \left| \left(1 + \frac{u}{n}\right)^n - 1 \right| \leq 2u, \quad \forall 0 \leq |u| \leq 1.$$

From this, we obviously have $|z - xy| \leq \epsilon |xy|$ by taking

$$(A.2) \quad \delta((x, y), \epsilon) = \frac{\epsilon}{6},$$

The output of any round-off machine if $x = 0$ or $y = 0$ is clearly 0, and hence the same value for ϵ of (A.2) suffices to satisfy the definition of computability.

EXAMPLE A.5. The same argument proves that the function $f : \mathbb{R}^\infty \rightarrow \mathbb{R}$ given by $f(x_1, \dots, x_n) = x_1 \cdots x_n$ is round-off computable (say, we compute first $x_1 x_2$ then $x_1 x_2 x_3$ and so on) with

$$(A.3) \quad \delta((x, y), \epsilon) = \frac{\epsilon}{4n - 2},$$

EXAMPLE A.6. A longer computation shows that the function $f : \{(x, y) \in \mathbb{R}^2 : x + y \neq 0\} \rightarrow \mathbb{R}$, $f(x, y) = x + y$ (again, we let it be zero in $\mathbb{R}^\infty \setminus \mathbb{R}^2$) is also round-off computable. It suffices to take

$$\delta((x, y), \epsilon) = \frac{\epsilon}{2 \max\left(1, \left|\frac{x}{x+y}\right|, \left|\frac{y}{x+y}\right|\right)}.$$

A more simple and still valid formula is

$$(A.4) \quad \delta((x, y), \epsilon) = \frac{|x + y|}{3\sqrt{2}\sqrt{x^2 + y^2}}\epsilon.$$

EXAMPLE A.7. Let us now see that $f(x) = x_1 + \dots + x_n$ is round-off computable in the set $\Omega = \{x \in \mathbb{R}^\infty : x_i \geq 0 \forall i\}$. Indeed, let $0 < \epsilon < 1$ and let us consider the most simple BSS machine which computes first $x_1 + x_2$, then adds x_3 and so on¹² A round-off machine with precision δ will produce, on input $x = (x_1, \dots, x_n)$, a number

$$x_1 \left(\prod_{k=1}^n (1 + \delta_1^{(k)}) \right) + x_2 \left(\prod_{k=1}^n (1 + \delta_2^{(k)}) \right) + \dots + x_n \left(\prod_{k=n-1}^n (1 + \delta_n^{(k)}) \right),$$

for some $\delta_i^{(k)}$ bounded in absolute value by δ . This follows from the fact that, in addition to the input error on each coordinate, x_1 and x_2 go through $n-1$ additions (which generate $n+1$ errors), x_3 goes through $n-2$ additions and so on. Note that

$$x_1(1 - \delta)^n \leq x_1 \left(\prod_{k=1}^n (1 + \delta_1^{(k)}) \right) \leq x_1(1 + \delta)^n.$$

Choosing $\delta = \alpha\epsilon/(2n)$, $0 < \alpha \leq 1$ and using (A.1) we conclude that

$$\left| x_1 \left(\prod_{k=1}^n (1 + \delta_1^{(k)}) \right) - x_1 \right| \leq \alpha\epsilon x_1,$$

and the same formula holds for x_2, \dots, x_n . The output of a round-off machine thus satisfies

$$\tilde{O}(x) = \sum_{i=1}^n x_i(1 + \alpha\epsilon_i), \quad 0 \leq |\epsilon_i| \leq \epsilon.$$

That is,

$$\left| \tilde{O}(x) - \sum_{i=1}^n x_i \right| = \sum_{i=1}^n x_i \alpha \epsilon_i \leq \sum_{i=1}^n x_i \alpha |\epsilon_i| \leq \alpha\epsilon \sum_{i=1}^n x_i,$$

proving that $f(x)$ is round-off computable in that set (just take $\alpha = 1$).

EXAMPLE A.8. Let us now see that $f(x) = x_1 + \dots + x_n$ is round-off computable in the set $\Omega = \{x \in \mathbb{R}^\infty : \sum x_i \neq 0\}$. We consider the BSS machine that first adds all the nonnegative numbers, call a the result, then adds all the negative numbers, call b the result, and then computes $a - b$. Let $0 < \epsilon < 1$. We note that from Example A.7 by choosing $\delta = \alpha\epsilon/(2n)$ (some $0 < \alpha \leq 1$) the round-off computation of the sum of positive (resp. negative) terms will be

$$\tilde{a} = a(1 + \alpha\epsilon_1), \quad \tilde{b} = b(1 + \alpha\epsilon_2), \quad \text{for some } 0 \leq |\epsilon_1|, |\epsilon_2| \leq \epsilon.$$

¹²This is not the algorithm of choice in practical programming but is sufficient for our purposes here.

From Example A.6, if we let

$$\alpha = \frac{|a + b|}{3\sqrt{2}\sqrt{a^2 + b^2}},$$

that is if we let

$$\delta(x, \epsilon) \leq \frac{|a + b|}{3\sqrt{2}\sqrt{a^2 + b^2}} \frac{\epsilon}{2n},$$

then $\tilde{O}(x) = \sum_i x_i$ up to relative precision ϵ . Using that $a^2 + b^2 \leq n \sum x_i^2$, we can also use the formula

$$(A.5) \quad \delta(x, \epsilon) = \frac{|\sum_{i=1}^n x_i|}{6\sqrt{2}n^{3/2}\sqrt{\sum_{i=1}^n x_i^2}} \epsilon.$$

EXAMPLE A.9. Combining examples A.5 and A.8 we see that the evaluation map of any multivariate polynomial $p(x_1, \dots, x_n)$ is round-off computable in the complement of its zero set (just compute first the monomials and then add all the results).

A.4. Ill-conditioned instances, condition number, posedness. Let us think of a function $f : \Omega \subseteq \mathbb{R}^\infty \rightarrow \mathbb{R}^\infty$ as the solution map associated with some problem to be solved. The condition number associated with f and x measures the first-order (relative) componentwise or normwise variations of $f(x)$ in terms of the first-order (relative) variations of x .

First assume that $f : \Omega \rightarrow \mathbb{R}$, that is the function is real-valued. We say that $x \in \bar{\Omega}$ (the topological closure of Ω) is well-conditioned when:

- Either $\|x\| \neq 0$, and f can be extended to a Lipschitz function defined in a neighborhood of x in $\bar{\Omega}$ with $|f(x)| \neq 0$. In that case we define the componentwise condition number by

$$\kappa_f(x) = \limsup_{x' \rightarrow x, x' \in \bar{\Omega}} \frac{\frac{|f(x') - f(x)|}{|f(x)|}}{\frac{\|x' - x\|}{\|x\|}},$$

- or f is constant in a neighborhood of x with $|f(x)| = 0$. In this later case we define the condition number by $\kappa_f(x) = 0$.

Otherwise, we say that $x \in \bar{\Omega}$ is ill-conditioned. The set of ill-conditioned instances is denoted by Σ_f , and for $x \in \Sigma_f$, we let $\kappa_f(x) = \infty$.

For a general $f : \Omega \rightarrow \mathbb{R}^\infty$, we define

$$\kappa_f(x) = \sup_j \kappa_{f_j}(x) \quad (\text{componentwise condition number})$$

that is the condition number of f is the supremum of the condition numbers of its coordinates. Sometimes it is more useful to consider the normwise condition number, that we denote by the same letter as the context should make clear which one is used on each problem:

$$\kappa_f(x) = \limsup_{x' \rightarrow x, x' \in \bar{\Omega}} \frac{\frac{\|f(x') - f(x)\|}{\|f(x)\|}}{\frac{\|x' - x\|}{\|x\|}} \quad (\text{normwise condition number}),$$

We define the posedness of a problem instance x with $\|x\| \neq 0$ as the distance to ill-posed problems:

$$\pi_f(x) = \frac{d(x, \Sigma_f)}{\|x\|}.$$

Here, $d(x, \Sigma_f) = \inf\{d(x, y) : y \in \Sigma_f\}$. The relation between condition number and posedness is an important but unclear problem. Following [33], we may expect a relation of the type

$$\pi_f(x) \approx \kappa_f(x)^{-1}$$

(condition number theorem) or at least inequalities like

$$C_1 \pi_f(x)^{\rho_1} \leq \kappa_f(x)^{-1} \leq C_w \pi_f(x)^{\rho_2}$$

for suitable positive constants C_i, ρ_i (cf. Lojasiewicz's inequality.) To get such a relation ill-posed problems should correspond to infinite condition numbers, but this is not always the case. Consider for example the decision problem: Is $x^2 + y^2 \leq \Pi$? The problem is well conditioned except on the circle $x^2 + y^2 = \Pi$, but the distance to this circle determines the precision we need in the computation.

Let $K_f(x) = \max(\kappa_f(x), \pi_f(x)^{-1})$.

EXAMPLE A.10. For $f(x) = x_1 \cdots x_n$ defined in \mathbb{R}^∞ , it is easy to see that

$$\kappa_f(x) = \sqrt{x_1^2 + \cdots + x_n^2} \sqrt{\frac{1}{x_1^2} + \cdots + \frac{1}{x_n^2}},$$

whenever $x_1, \dots, x_n \neq 0$. If $x_i = 0$ for any i then $\kappa_f(x) = \infty$.

On the other hand,

$$\pi_f(x) = \frac{\min(|x_1|, \dots, |x_n|)}{\sqrt{x_1^2 + \cdots + x_n^2}}.$$

Thus, we have

$$\kappa_f(x) \leq \sqrt{x_1^2 + \cdots + x_n^2} \sqrt{\frac{n}{\min(|x_1|, \dots, |x_n|)^2}} = \sqrt{n} \pi_f(x)^{-1},$$

and

$$\kappa_f(x) \geq \sqrt{x_1^2 + \cdots + x_n^2} \sqrt{\frac{1}{\min(|x_1|, \dots, |x_n|)^2}} = \pi_f(x)^{-1}.$$

Namely,

$$\pi_f(x)^{-1} \leq \kappa_f(x) \leq \sqrt{n} \pi_f(x)^{-1}.$$

EXAMPLE A.11. For $f(x) = x_1 + \cdots + x_n$ defined in $\Omega = \{x \in \mathbb{R}^\infty : \sum x_i \neq 0\}$, we have:

- For $x \in \Omega$, a simple computation shows that

$$\kappa_f(x) = \frac{\sqrt{n} \sqrt{\sum x_i^2}}{|\sum x_i|}.$$

- For $x \in \partial\Omega$, that is $\sum x_i = 0$, we have $\kappa_f(x) = \infty$.

Thus, we have

$$\pi_f(x) = \frac{d(x, \{x : \sum x_i = 0\})}{\sqrt{\sum x_i^2}} = \frac{|\sum x_i|}{\sqrt{n} \sqrt{\sum x_i^2}} = \kappa_f(x, y)^{-1}.$$

Namely,

$$(A.6) \quad K_f(x) = \frac{\sqrt{n} \sqrt{\sum x_i^2}}{|\sum x_i|}.$$

A.5. Single, multiple precision. Let f be a round-off computable function, and let M be a BSS machine satisfying the definition of round-off computability above. This computation is single precision when for every $0 < \epsilon < 1$ there is a $\delta = \delta(\epsilon)$ such that any round-off machine (M, δ) attains relative precision ϵ for any input $x \in \Omega$, and such that

$$(A.7) \quad \delta \geq \frac{c_0 \epsilon}{K_f(x)^{c_2} \dim(x)^{c_3}}$$

for some positive constants c_0, c_2, c_3 . This computation is multiple precision when there exists δ such that

$$(A.8) \quad \delta \geq \frac{c_0 \epsilon^{c_1}}{K_f(x)^{c_2} \dim(x)^{c_3}},$$

for some $c_1 > 1$. We say that the computation is strictly multiple precision when it is multiple precision but not single precision.

EXAMPLE A.12. The inductive, naive algorithm for computing the round-off computable function $f(x_1, \dots, x_n) = x_1 \cdots x_n$ defined in \mathbb{R}^∞ is single precision, from (A.3). The algorithm given in Example A.8 for computing the round-off computable function $f(x) = x_1 + \cdots + x_n$ defined in $\{x \in \mathbb{R}^\infty : \sum x_i \neq 0\}$ is single precision from (A.5) and (A.6).

A.6. Size of an input. In many practical problems, we want to specify an output precision ϵ . From our definition of round-off computable function, given $x \in \Omega$ and $0 < \epsilon < 1$ some $\delta(x, \epsilon)$ will exist guaranteeing the desired precision, although it may be very hard to compute this δ in some cases. Moreover, from (A.8), the number $K_f(x)$ will in general play a role in the value of $\delta(x, \epsilon)$ needed for any machine solving the problem. This dependence suggests that maybe the input should be considered as (x, ϵ) and not just as x . These thoughts justify our definition of the size of an input, which includes a term related to ϵ and another related to $K_f(x)$:

$$(A.9) \quad \dim(x) + |\log \epsilon| + \log(K_f(x) + 1).$$

A.7. Cost of a computation. The cost of a computation on a round-off machine (M, δ) which outputs \tilde{y} on input x is

$$T(x, \delta) \cdot \left(\max_i \dim(y^{(i)}) + |\log \delta| \right),$$

where $T(x, \delta)$ is the time for the computation to halt and

$$x = y^{(0)}, \dots, y^{T(x, \delta)} = \tilde{y}$$

are the different vectors computed by (M, δ) on input x .

We say that a function $f : \Omega \rightarrow \mathbb{R}^\infty$ is polynomial cost computable if there exists a BSS machine M such that for every $x \in \Omega$ and $0 < \epsilon < 1$ there exists $\delta(x, \epsilon)$ such that any round-off machine $(M, \delta(x, \epsilon))$ computes \tilde{y} which equals $f(x)$ to relative error ϵ , with cost polynomially bounded by the input size (A.9).

The most important cases of polynomial cost computability will be in the cases where we restrict the space of functions to single (multiple) precision functions, for example in the case of single precision to the definition of polynomial cost we add the restriction that $\delta(x, \epsilon)$ must satisfy (A.7). These two possibilities (single or multiple precision) will give us two theories, both of which deserve to be worked out.

Now that we have the notion of polynomial cost the classes P and NP may be defined and the problem: Does $P = NP$? stated.

References

- [1] D. Armentano. Complexity of path-following methods for the eigenvalue problem . To appear.
- [2] D. Armentano. PH. D. Thesis. Universidad de la República, Uruguay, and Université Paul Sabatier, France.
- [3] Diego Armentano, Carlos Beltrán, and Michael Shub, *Minimizing the discrete logarithmic energy on the sphere: the role of random polynomials*, Trans. Amer. Math. Soc. **363** (2011), no. 6, 2955–2965, DOI 10.1090/S0002-9947-2011-05243-8. MR2775794 (2012f:31009)
- [4] D. Armentano, M. Shub. Smale’s fundamental theorem of algebra reconsidered. To appear in Foundations of Computational Mathematics. DOI: 10.1007/s10208-013-9155-y.
- [5] Steve Batterson, *Convergence of the Francis shifted QR algorithm on normal matrices*, Linear Algebra Appl. **207** (1994), 181–195, DOI 10.1016/0024-3795(94)90010-8. MR1283957 (95h:65028)
- [6] Steve Batterson and John Smillie, *Rayleigh quotient iteration fails for nonsymmetric matrices*, Appl. Math. Lett. **2** (1989), no. 1, 19–20, DOI 10.1016/0893-9659(89)90107-9. MR989851 (90b:65086)
- [7] D. J. Bates, J. D. Hauenstein, A. J. Sommese, and C. W. Wampler. Bertini: software for numerical algebraic geometry. Available at <http://www.nd.edu/~sommese/bertini>.
- [8] Carlos Beltrán, *A continuation method to solve polynomial systems and its complexity*, Numer. Math. **117** (2011), no. 1, 89–113, DOI 10.1007/s00211-010-0334-3. MR2754220 (2011m:65102)
- [9] C. Beltrán. The state of the art in Smale’s 7–th problem. In Foundations of Computational Mathematics, Budapest 2011. London Mathematical Society. Lecture notes series 403. F. Cucker, T. Krick, A. Pinkus, A. Szanto editors.
- [10] Carlos Beltrán, Jean-Pierre Dedieu, Gregorio Malajovich, and Mike Shub, *Convexity properties of the condition number*, SIAM J. Matrix Anal. Appl. **31** (2009), no. 3, 1491–1506, DOI 10.1137/080718681. MR2587788 (2011c:65071)
- [11] Carlos Beltrán, Jean-Pierre Dedieu, Gregorio Malajovich, and Mike Shub, *Convexity properties of the condition number II*, SIAM J. Matrix Anal. Appl. **33** (2012), no. 3, 905–939, DOI 10.1137/100808885. MR3023457
- [12] Carlos Beltrán and Anton Leykin, *Certified numerical homotopy tracking*, Exp. Math. **21** (2012), no. 1, 69–83, DOI 10.1080/10586458.2011.606184. MR2904909
- [13] Carlos Beltrán and Anton Leykin, *Robust Certified Numerical Homotopy Tracking*, Found. Comput. Math. **13** (2013), no. 2, 253–295, DOI 10.1007/s10208-013-9143-2. MR3032682
- [14] Carlos Beltrán and Luis Miguel Pardo, *On Smale’s 17th problem: a probabilistic positive solution*, Found. Comput. Math. **8** (2008), no. 1, 1–43, DOI 10.1007/s10208-005-0211-0. MR2403529 (2009h:65082)
- [15] Carlos Beltrán and Luis Miguel Pardo, *Smale’s 17th problem: average polynomial time to compute affine and projective solutions*, J. Amer. Math. Soc. **22** (2009), no. 2, 363–385, DOI 10.1090/S0894-0347-08-00630-9. MR2476778 (2009m:90147)
- [16] Carlos Beltrán and Luis Miguel Pardo, *Fast linear homotopy to find approximate zeros of polynomial systems*, Found. Comput. Math. **11** (2011), no. 1, 95–129, DOI 10.1007/s10208-010-9078-9. MR2754191 (2011m:65111)
- [17] Carlos Beltrán and Michael Shub, *Complexity of Bezout’s theorem. VII. Distance estimates in the condition metric*, Found. Comput. Math. **9** (2009), no. 2, 179–195, DOI 10.1007/s10208-007-9018-5. MR2496559 (2010f:65100)
- [18] Carlos Beltrán and Michael Shub, *A note on the finite variance of the averaging function for polynomial system solving*, Found. Comput. Math. **10** (2010), no. 1, 115–125, DOI 10.1007/s10208-009-9054-4. MR2591841 (2011b:65075)
- [19] Carlos Beltrán and Michael Shub, *On the geometry and topology of the solution variety for polynomial system solving*, Found. Comput. Math. **12** (2012), no. 6, 719–763, DOI 10.1007/s10208-012-9134-8. MR2989472
- [20] Lenore Blum, Felipe Cucker, Michael Shub, and Steve Smale, *Complexity and real computation*, Springer-Verlag, New York, 1998. With a foreword by Richard M. Karp. MR1479636 (99a:68070)

- [21] Lenore Blum, Mike Shub, and Steve Smale, *On a theory of computation and complexity over the real numbers: NP-completeness, recursive functions and universal machines*, Bull. Amer. Math. Soc. (N.S.) **21** (1989), no. 1, 1–46, DOI 10.1090/S0273-0979-1989-15750-9. MR974426 (90a:68022)
- [22] Paola Boito and Jean-Pierre Dedieu, *The condition metric in the space of rectangular full rank matrices*, SIAM J. Matrix Anal. Appl. **31** (2010), no. 5, 2580–2602, DOI 10.1137/08073874X. MR2740622 (2012e:65078)
- [23] Allan Borodin and Ian Munro, *The computational complexity of algebraic and numeric problems*, American Elsevier Publishing Co., Inc., New York-London-Amsterdam, 1975. Elsevier Computer Science Library; Theory of Computation Series, No. 1. MR0468309 (57 #8145)
- [24] J. S. Brauchart, *Optimal logarithmic energy points on the unit sphere*, Math. Comp. **77** (2008), no. 263, 1599–1613, DOI 10.1090/S0025-5718-08-02085-1. MR2398782 (2010e:31004)
- [25] M. Braverman. On the complexity of real functions, FOCS 2005.
- [26] Mark Braverman and Stephen Cook, *Computing over the reals: foundations for scientific computing*, Notices Amer. Math. Soc. **53** (2006), no. 3, 318–329. MR2208383 (2006m:68019)
- [27] Peter Bürgisser and Felipe Cucker, *On a problem posed by Steve Smale*, Ann. of Math. (2) **174** (2011), no. 3, 1785–1836, DOI 10.4007/annals.2011.174.3.8. MR2846491
- [28] P. Bürgisser and F. Cucker. *Condition: The Geometry of Numerical Algorithms*. Grundlehren der mathematischen Wissenschaften, 349. ISBN-10:3642388957 — ISBN-13: 978-3642388958.
- [29] Felipe Cucker and Steve Smale, *Complexity estimates depending on condition and round-off error*, J. ACM **46** (1999), no. 1, 113–184, DOI 10.1145/300515.300519. MR1692497 (2000f:68040)
- [30] F. Cucker and J.-P. Dedieu, *Decision problems and round-off machines*, Theory Comput. Syst. **34** (2001), no. 5, 433–452. MR1862890 (2002h:68050)
- [31] Jean-Pierre Dedieu, Gregorio Malajovich, and Michael Shub, *Adaptive step-size selection for homotopy methods to solve polynomial equations*, IMA J. Numer. Anal. **33** (2013), no. 1, 1–29, DOI 10.1093/imanum/drs007. MR3020948
- [32] Jean-Pierre Dedieu and Mike Shub, *Multihomogeneous Newton methods*, Math. Comp. **69** (2000), no. 231, 1071–1098 (electronic), DOI 10.1090/S0025-5718-99-01114-X. MR1752092 (2000m:65072)
- [33] James Weldon Demmel, *On condition numbers and the distance to the nearest ill-posed problem*, Numer. Math. **51** (1987), no. 3, 251–289, DOI 10.1007/BF01400115. MR895087 (88i:15014)
- [34] A. Dubickas, *On the maximal product of distances between points on a sphere*, Liet. Mat. Rink. **36** (1996), no. 3, 303–312, DOI 10.1007/BF02986850 (English, with English and Lithuanian summaries); English transl., Lithuanian Math. J. **36** (1996), no. 3, 241–248 (1997). MR1455810 (98e:52015)
- [35] C. Eckart and G. Young. The approximation of one matrix by another of lower rank, Psychometrika 1 (1936), 211–218.
- [36] Misha Gromov, *Metric structures for Riemannian and non-Riemannian spaces*, Progress in Mathematics, vol. 152, Birkhäuser Boston Inc., Boston, MA, 1999. Based on the 1981 French original [MR0682063 (85e:53051)]; With appendices by M. Katz, P. Pansu and S. Semmes; Translated from the French by Sean Michael Bates. MR1699320 (2000d:53065)
- [37] G.H. Hardy, J.E. Littlewood, G. Pólya. Inequalities, Cambridge University Press, 1934.
- [38] Nicholas J. Higham, *Accuracy and stability of numerical algorithms*, 2nd ed., Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2002. MR1927606 (2003g:65064)
- [39] M.H. Kim. Computational complexity of the Euler type algorithms for the roots of complex polynomials. PhD thesis, The City University of New York, 1985.
- [40] Myong-Hi Kim, *On approximate zeros and rootfinding algorithms for a complex polynomial*, Math. Comp. **51** (1988), no. 184, 707–719, DOI 10.2307/2008771. MR958638 (90f:65073)
- [41] T. L. Lee, T. Y. Li, and C. H. Tsai. Hom4ps-2.0: A software package for solving polynomial systems by the polyhedral homotopy continuation method. Available at http://hom4ps.math.msu.edu/HOM4PS_soft.htm.
- [42] Anton Leykin, *Numerical algebraic geometry*, J. Softw. Algebra Geom. **3** (2011), 5–10. MR2881262
- [43] Gregorio Malajovich, *Nonlinear equations*, Publicações Matemáticas do IMPA. [IMPA Mathematical Publications], Instituto Nacional de Matemática Pura e Aplicada (IMPA), Rio de

- Janeiro, 2011. With an appendix by Carlos Beltrán, Jean-Pierre Dedieu, Luis Miguel Pardo and Mike Shub; 28° Colóquio Brasileiro de Matemática. [28th Brazilian Mathematics Colloquium]. MR2798351 (2012j:65148)
- [44] C.W. Pfrang, P. Deift and G. Menon. How long does it take to compute the eigenvalues of a random symmetric matrix? arXiv 1203.4635.
- [45] E. A. Rakhmanov, E. B. Saff, and Y. M. Zhou, *Minimal discrete energy on the sphere*, Math. Res. Lett. **1** (1994), no. 6, 647–662. MR1306011 (96e:78011)
- [46] James Renegar, *On the cost of approximating all roots of a complex polynomial*, Math. Programming **32** (1985), no. 3, 319–336, DOI 10.1007/BF01582052. MR796429 (87a:65083)
- [47] James Renegar, *On the worst-case arithmetic complexity of approximating zeros of polynomials*, J. Complexity **3** (1987), no. 2, 90–113, DOI 10.1016/0885-064X(87)90022-7. MR907192 (89a:68107)
- [48] James Renegar, *On the worst-case arithmetic complexity of approximating zeros of systems of polynomials*, SIAM J. Comput. **18** (1989), no. 2, 350–370, DOI 10.1137/0218024. MR986672 (90j:68021)
- [49] James Renegar, *Incorporating condition measures into the complexity theory of linear programming*, SIAM J. Optim. **5** (1995), no. 3, 506–524, DOI 10.1137/0805026. MR1344668 (96c:90048)
- [50] Michael Shub, *Some remarks on Bezout's theorem and complexity theory*, From Topology to Computation: Proceedings of the Smalefest (Berkeley, CA, 1990), Springer, New York, 1993, pp. 443–455. MR1246139 (95a:14002)
- [51] Michael Shub, *Complexity of Bezout's theorem. VI. Geodesics in the condition (number) metric*, Found. Comput. Math. **9** (2009), no. 2, 171–178, DOI 10.1007/s10208-007-9017-6. MR2496558 (2010f:65103)
- [52] Michael Shub and Steve Smale, *Complexity of Bézout's theorem. I. Geometric aspects*, J. Amer. Math. Soc. **6** (1993), no. 2, 459–501, DOI 10.2307/2152805. MR1175980 (93k:65045)
- [53] M. Shub and S. Smale, *Complexity of Bezout's theorem. II. Volumes and probabilities*, Computational algebraic geometry (Nice, 1992), Progr. Math., vol. 109, Birkhäuser Boston, Boston, MA, 1993, pp. 267–285. MR1230872 (94m:68086)
- [54] Michael Shub and Steve Smale, *Complexity of Bezout's theorem. III. Condition number and packing*, J. Complexity **9** (1993), no. 1, 4–14, DOI 10.1006/jcom.1993.1002. Festschrift for Joseph F. Traub, Part I. MR1213484 (94g:65152)
- [55] Michael Shub and Steve Smale, *Complexity of Bezout's theorem. IV. Probability of success; extensions*, SIAM J. Numer. Anal. **33** (1996), no. 1, 128–148, DOI 10.1137/0733008. MR1377247 (97k:65310)
- [56] M. Shub and S. Smale, *Complexity of Bezout's theorem. V. Polynomial time*, Theoret. Comput. Sci. **133** (1994), no. 1, 141–164, DOI 10.1016/0304-3975(94)90122-8. Selected papers of the Workshop on Continuous Algorithms and Complexity (Barcelona, 1993). MR1294430 (96d:65091)
- [57] Steve Smale, *The fundamental theorem of algebra and complexity theory*, Bull. Amer. Math. Soc. (N.S.) **4** (1981), no. 1, 1–36, DOI 10.1090/S0273-0979-1981-14858-8. MR590817 (83i:65044)
- [58] Steve Smale, *Newton's method estimates from data at one point*, computational mathematics (Laramie, Wyo., 1985), Springer, New York, 1986, pp. 185–196. MR870648 (88e:65076)
- [59] S. Smale. The fundamental theorem of algebra and complexity theory, SIAM Rev. **32** (1990), no. 2, 211–220.
- [60] Steve Smale, *Mathematical problems for the next century*, Mathematics: frontiers and perspectives, Amer. Math. Soc., Providence, RI, 2000, pp. 271–294. MR1754783 (2001i:00003)
- [61] G. W. Stewart and Ji Guang Sun, *Matrix perturbation theory*, Computer Science and Scientific Computing, Academic Press Inc., Boston, MA, 1990. MR1061154 (92a:65017)
- [62] Lloyd N. Trefethen and David Bau III, *Numerical linear algebra*, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1997. MR1444820 (98k:65002)
- [63] J. Verschelde. Algorithm 795: PHCpack: A general-purpose solver for polynomial systems by homotopy continuation. ACM Trans. Math. Softw., **25** (1999), no. 2, 251–276. Available at <http://www.math.uic.edu/~jan>.
- [64] Gerold Wagner, *On the product of distances to a point set on a sphere*, J. Austral. Math. Soc. Ser. A **47** (1989), no. 3, 466–482. MR1018975 (90j:11080)

- [65] Qi Zhong, *Energies of zeros of random sections on Riemann surfaces*, *Indiana Univ. Math. J.* **57** (2008), no. 4, 1753–1780, DOI 10.1512/iumj.2008.57.3329. MR2440880 (2009k:58051)
- [66] Steve Zelditch and Qi Zhong, *Addendum to “Energies of zeros of random sections on Riemann surfaces”*. *Indiana Univ. Math. J.* **57** (2008), No. 4, 1753–1780 [MR 2440880], *Indiana Univ. Math. J.* **59** (2010), no. 6, 2001–2005, DOI 10.1512/iumj.2010.59.59073. MR2919745

DEPTO. DE MATEMÁTICAS, ESTADÍSTICA Y COMPUTACIÓN, UNIVERSIDAD DE CANTABRIA, SANTANDER, SPAIN.

E-mail address: `carlos.beltran@unican.es`

CONICET, IMAS, UNIVERSIDAD DE BUENOS AIRES, ARGENTINA AND CUNY GRADUATE SCHOOL, NEW YORK, NY, USA.

E-mail address: `shub.michael@gmail.com`