

On Approximate Triangular Decompositions of Zero Dimensional Systems

Marc Moreno Maza, Greg Reid, Robin Scott and Wenyuan Wu

Abstract. Triangular decompositions for systems of polynomial equations with n variables, with exact coefficients are well-developed theoretically and in terms of implemented algorithms in computer algebra systems. However there is much less research about triangular decompositions for systems with approximate coefficients.

In this paper we discuss the zero dimensional case, of systems having finitely many roots. Our methods depend on having approximations for all the roots, and these are provided by the homotopy continuation methods of Sommese, Verschelde and Wampler. We introduce approximate equiprojectable decompositions for such systems, which represent a generalization of the recently developed analogous concept for exact systems. We demonstrate experimentally the favourable computational features of this new approach, and give a statistical analysis of its error.

Keywords. Symbolic-numeric computations, Triangular decompositions, Dimension zero, Polynomial system solving.

1. Introduction

Ritt initiated the algebraic study of differential polynomial systems through characteristic sets [22]. Their modern study was revitalized by the work of Wu. In [32], he adapted the work of Ritt for solving algebraic systems: he showed that the zero set of such a system could be decomposed as finitely many characteristic sets, leading to the notion of a triangular decomposition of an algebraic variety. Considerable developments have followed by many authors; among them: Chou [5], Dahan et al. [7], Gao et al. [11], Kalkbrener [13], Lazard [14], Moreno Maza [19], Schost [23], Wang [31], and others. These works have led to efficient algorithms for triangular decomposition of an algebraic variety given by an exact input polynomial system.

This work is supported by NSERC, MITACS and Maplesoft, Canada.

Often, in applications we are interested in producing a useful triangular form where some of the variables are functions of others. Such systems frequently have approximate coefficients that are inferred from experimental data. This means that the stability, or sensitivity to coefficient changes, of such triangular decompositions is a concern. While considerable progress in both theoretical and algorithmic aspects has been made for exact input polynomial systems, much less is known about generalizations of these methods to input systems which are approximate.

In this paper, we present some initial results in this direction, for the case of an algebraic variety V over \mathbb{C} . We rely on the methods of Sommese, Verschelde, and Wampler [25, 30, 18, 26] which use Homotopy continuation, to determine so-called generic points on the components of the numerical decomposition of V . We are interested in the set V_0 of the isolated points of V (the 0 dimensional case). Each point of V_0 , and more generally every irreducible component of V , is trivially a triangular set, although not generally rationally constructible from rational input. This is in contrast to the usual forms of exact triangular decomposition, which are modeled on equi-dimensional decomposition over \mathbb{Q} rather than irreducible decomposition over \mathbb{C} .

Following [6, 7], we consider the equiprojectable decomposition of V_0 . Then, we use the interpolation formulas of Dahan and Schost [8] for computing an approximate triangular set for each equiprojectable component of V_0 , leading to an approximate triangular decomposition of V_0 in Section 3.

We provide a stability analysis of the interpolation formulas of Dahan and Schost in Section 4. One of our main tools is Lindeberg's theorem [24] that is described in the Appendix. In Sections 5 and 6, we report on experiments that illustrate the efficiency of our approach and support the accuracy of our stability analysis.

In [21], we study the simplest class of positive dimensional systems: linear homogeneous systems. Our aim in that article is to explore local structure of nonlinear problems with linearized approximate triangular decompositions. The combination of the two approaches allows us to form an accessible bridge to the study of the fully non-linear case which we will describe in a forthcoming paper.

2. Triangular decompositions

A triangular decomposition of a zero-dimensional algebraic variety V is a family of polynomial sets, called triangular sets, that describe symbolically the points of V [14]. Triangular decompositions extend to algebraic varieties of arbitrary dimension, see for instance [13, 19]. In [8] it is shown that the height of a coefficient in a triangular set T can be bounded by the height of the variety represented by T . Combined with the notion of *equiprojectable decomposition* introduced in [6], this motivated the work of [7], in which the authors obtained a very efficient method for computing triangular decompositions of zero-dimensional varieties over \mathbb{Q} given by an input polynomial system with exact coefficients.

On top of these good computational properties, triangular sets and triangular decompositions have natural geometrical interpretations. In Section 3, we will rely on these properties to introduce a notion of an *approximate triangular decomposition* of a zero-dimensional variety given by approximate coordinates of its points. In the present section, we recall some results for triangular decompositions in the exact case and refer to [8, 6, 7] for more details. For the reader's convenience, we sketch the proof of Propositions 2.4 and 2.5, which play a central role in this paper. See [8] for their complete proofs.

Let \mathbb{K} be a perfect field, let \mathbb{L} be an algebraic closure of \mathbb{K} and let $X_1 \prec \dots \prec X_n$ be $n \geq 1$ ordered variables.

Definition 2.1. A set $T = \{T_1, \dots, T_n\}$ of n polynomials in $\mathbb{K}[X_1, \dots, X_n]$ is a *triangular set* if the ideal $\langle T \rangle$ generated by T is radical and if for all $1 \leq i \leq n$ the polynomial T_i is not constant, the greatest variable occurring in T_i is X_i , and its leading coefficient w.r.t. X_i is invertible modulo the ideal $\langle T_1, \dots, T_{i-1} \rangle$. The triangular set T is *normalized* if for all $1 \leq i \leq n$ the leading coefficient of T_i w.r.t. X_i is one.

Clearly, a triangular set generates a zero-dimensional ideal and a normalized triangular set is a reduced lexicographical Gröbner basis. In [14], it is shown that every maximal ideal of $\mathbb{K}[X_1, \dots, X_n]$ can be generated by a triangular set. Hence, a natural question is to characterize the zero-dimensional varieties over \mathbb{K} , that can be generated by a triangular set. The answer is given by [2]. We report on it here by means of Definition 2.2 and Theorem 2.3, after introducing some notation.

Let i and j be integers such that $1 \leq i \leq j \leq n$. We denote by $A^i(\mathbb{L})$ the affine space of dimension i over \mathbb{L} . For $V \subseteq A^n(\mathbb{L})$ we denote by $\mathcal{I}(V)$ the ideal of $\mathbb{K}[X_1, \dots, X_n]$ composed by the polynomials which vanish on V . For $F \subseteq \mathbb{K}[X_1, \dots, X_n]$ we denote by $V(F)$ the set of the points of $A^n(\mathbb{L})$ where every element of F vanishes. Finally, we denote by π_i^j the natural projection map from $A^j(\mathbb{L})$ to $A^i(\mathbb{L})$, which sends (X_1, \dots, X_j) to (X_1, \dots, X_i) .

Definition 2.2. A zero-dimensional variety $V \subseteq A^j(\mathbb{L})$ over \mathbb{K} is said to be

- (1) *equiprojectable on* $V_i = \pi_i^j(V)$, its projection onto $A^i(\mathbb{L})$, if there exists an integer c such that for every $M \in V_i$ the cardinality of $(\pi_i^j)^{-1}(M) \cap V$ is c .
- (2) *equiprojectable* if V is equiprojectable on V_1, \dots, V_{j-1} .

Theorem 2.3. *A zero-dimensional variety $V \subseteq A^j(\mathbb{L})$ over \mathbb{K} is equiprojectable if and only if there exists a triangular set T of $\mathbb{K}[X_1, \dots, X_j]$ such that T generates $\mathcal{I}(V)$.*

Given an equiprojectable variety $V \subseteq A^n(\mathbb{L})$ the normalized triangular set T generating $\mathcal{I}(V)$ can be constructed as follows from the coordinates of the points of V (see [8] for details). Let \mathbf{K} be a field such that $\mathbb{K} \subseteq \mathbf{K} \subseteq \mathbb{L}$ and such that every point of V has its coordinates in \mathbf{K} . We define $V_i = \pi_i^n(V)$. Let $1 \leq \ell < n$. Following [8], we describe how to interpolate $T_{\ell+1}$ from the coordinates (in \mathbf{K}) of

the points of $V_{\ell+1}$. Let $\alpha = (\alpha_1, \dots, \alpha_\ell) \in V_\ell$. Define:

$$\begin{aligned} V_\alpha^1 &= \{\beta = (\beta_1, \dots, \beta_\ell, \beta_{\ell+1}) \in V_{\ell+1} \mid \beta_1 \neq \alpha_1\}, \\ V_\alpha^2 &= \{\beta = (\alpha_1, \beta_2, \dots, \beta_\ell, \beta_{\ell+1}) \in V_{\ell+1} \mid \beta_2 \neq \alpha_2\}, \\ V_\alpha^3 &= \{\beta = (\alpha_1, \alpha_2, \beta_3, \dots, \beta_\ell, \beta_{\ell+1}) \in V_{\ell+1} \mid \beta_3 \neq \alpha_3\}, \\ &\dots \quad \dots \quad \dots \\ V_\alpha^\ell &= \{\beta = (\alpha_1, \dots, \alpha_{\ell-1}, \beta_\ell, \beta_{\ell+1}) \in V_{\ell+1} \mid \beta_\ell \neq \alpha_\ell\}, \\ V_\alpha^{\ell+1} &= \{\beta = (\alpha_1, \dots, \alpha_\ell, \beta_{\ell+1}) \in V_{\ell+1}\}. \end{aligned} \quad (2.1)$$

The sets $V_\alpha^1, V_\alpha^2, V_\alpha^3, \dots, V_\alpha^\ell, V_\alpha^{\ell+1}$ partition $V_{\ell+1}$. We consider also the projections:

$$\begin{aligned} v_\alpha^1 &= \pi_1^{\ell+1}(V_\alpha^1) = \{(\beta_1) \in V_1 \mid \beta_1 \neq \alpha_1\}, \\ v_\alpha^2 &= \pi_2^{\ell+1}(V_\alpha^2) = \{(\alpha_1, \beta_2) \in V_2 \mid \beta_2 \neq \alpha_2\}, \\ v_\alpha^3 &= \pi_3^{\ell+1}(V_\alpha^3) = \{(\alpha_1, \alpha_2, \beta_3) \in V_3 \mid \beta_3 \neq \alpha_3\}, \\ &\dots \quad \dots \quad \dots \quad \dots \quad \dots \\ v_\alpha^\ell &= \pi_\ell^{\ell+1}(V_\alpha^\ell) = \{(\alpha_1, \dots, \alpha_{\ell-1}, \beta_\ell) \in V_\ell \mid \beta_\ell \neq \alpha_\ell\} \end{aligned} \quad (2.2)$$

For $1 \leq i \leq \ell + 1$, we define

$$T_{\alpha,i} = T_i(\alpha_1, \dots, \alpha_{i-1}, X_i) \quad \text{and} \quad e_{\alpha,i} = \prod_{\beta \in v_\alpha^i} (X_i - \beta_i). \quad (2.3)$$

Observe that for $1 \leq i \leq \ell + 1$ we have $T_{\alpha,i} \in \mathbf{K}[X_i]$ and $e_{\alpha,i} \in \mathbf{K}[X_i]$. Finally, we define

$$E_\alpha = \prod_{1 \leq i \leq \ell} e_{\alpha,i} \quad (2.4)$$

and note that $E_\alpha \in \mathbf{K}[X_1, \dots, X_\ell]$ holds.

Proposition 2.4. *For $1 \leq i \leq \ell$ we have*

$$T_{\alpha,i} = \prod_{(\alpha_1, \dots, \alpha_{i-1}, \beta_i) \in V_i} (X_i - \beta_i) = e_{\alpha,i} (X_i - \alpha_i), \quad (2.5)$$

$$T_{\alpha,\ell+1} = \prod_{\beta \in V_\alpha^{\ell+1}} (X_{\ell+1} - \beta_{\ell+1}), \quad (2.6)$$

$$T_{\ell+1} = \sum_{\alpha \in V_\ell} \frac{E_\alpha T_{\alpha,\ell+1}}{E_\alpha(\alpha)}. \quad (2.7)$$

PROOF. Relations (2.5) and (2.6) follow easily (2.1), (2.2) and (2.3). In order to prove (2.7) we observe that:

$$(\forall \beta \in V_\ell) \quad E_\alpha(\beta) = 0 \iff \beta \neq \alpha. \quad (2.8)$$

Indeed, for $1 \leq i \leq \ell$, we have $e_{\alpha,i}(\alpha) \neq 0$ leading to $E_\alpha(\alpha) \neq 0$. Now let $\beta \in V_\ell$ with $\beta \neq \alpha$. Then, there exists $i \leq \ell$ such that

$$(\pi_i^\ell)^{-1}(\beta) \in v_\alpha^i.$$

Hence, for this index i we have $e_{\alpha,i}(\beta) = 0$, which proves (2.8). From there, establishing (2.7) is routine. \square

In [8], another triangular set N is obtained from the coordinates of the points of V , see Proposition 2.5. The authors show that it has much smaller coefficients than the normalized triangular set given by the formulas of Proposition 2.4. We will be generalizing this second triangular set to the approximate case.

Proposition 2.5 (Interpolation formulas). *Let $D_1 = 1$ and $\tau_1 = N_1 = T_1$. For $2 \leq \ell \leq n$, define*

$$D_\ell = \prod_{1 \leq i \leq \ell-1} \frac{\partial T_i}{\partial X_i} \pmod{\langle T_1, \dots, T_{\ell-1} \rangle} \quad (2.9)$$

and

$$N_\ell = D_\ell T_\ell \pmod{\langle T_1, \dots, T_{\ell-1} \rangle}. \quad (2.10)$$

Then, for $1 \leq i \leq \ell$ we have

$$N_{\ell+1} = \sum_{\alpha \in V_\ell} E_\alpha T_{\alpha, \ell+1}. \quad (2.11)$$

PROOF. Indeed, for $1 \leq i \leq \ell$, we have

$$T_{\alpha, i} = e_{\alpha, i} (X_i - \alpha_i) \in \mathbf{K}[X_i]$$

leading to

$$\begin{aligned} \frac{\partial T}{\partial X_i}(\alpha) &= T'_{\alpha, i}(\alpha) \\ &= e'_{\alpha, i}(\alpha) (\alpha_i - \alpha_i) + e_{\alpha, i}(\alpha) \\ &= e_{\alpha, i}(\alpha). \end{aligned}$$

By definition, we have

$$N_{\ell+1} = \left(\prod_{1 \leq i \leq \ell} \frac{\partial T}{\partial X_i} \right) T_{\ell+1} \pmod{\langle T_1, \dots, T_\ell \rangle}.$$

Hence, we have

$$\begin{aligned} N_{\ell+1}(\alpha) &= \left(\prod_{1 \leq i \leq \ell} \frac{\partial T}{\partial X_i}(\alpha) \right) T_{\ell+1}(\alpha) \\ &= \left(\prod_{1 \leq i \leq \ell} e_{\alpha, i}(\alpha) \right) T_{\ell+1}(\alpha) \\ &= E_\alpha(\alpha) T_{\ell+1}(\alpha) \end{aligned}$$

where $T_{\ell+1}(\alpha) = T_{\alpha, \ell+1}$ holds. Finally we obtain

$$\begin{aligned} N_{\ell+1} &= \sum_{\alpha \in V_\ell} \frac{E_\alpha N_{\ell+1}(\alpha)}{E_\alpha(\alpha)} \\ &= \sum_{\alpha \in V_\ell} E_\alpha T_{\ell+1}(\alpha). \end{aligned}$$

□

Clearly, not all zero-dimensional varieties over \mathbb{Q} are equiprojectable. Consider, for instance, with $n = 2$ the variety consisting of the three points A , B , C with respective coordinates $(1, 0)$, $(0, 0)$ and $(0, 1)$. However, we do have the following result, see for instance [14].

Proposition 2.6. *For every radical ideal \mathcal{I} of $\mathbb{K}[X_1, \dots, X_n]$ there exists finitely many triangular sets T^1, \dots, T^e such that \mathcal{I} is the intersection of the ideals $\langle T^1 \rangle, \dots, \langle T^e \rangle$. If, in addition, the ideals $\langle T^1 \rangle, \dots, \langle T^e \rangle$ are pairwise relatively prime, then the set $\{T^1, \dots, T^e\}$ is called a triangular decomposition of the ideal \mathcal{I} .*

Triangular decompositions of algebraic varieties (with arbitrary dimension) are discussed in depth in [19] together with an algorithm for computing them, which is implemented in [16]. Observe that a radical ideal may admit several triangular decompositions. For instance, there are four different triangular decompositions for the ideal $\mathcal{I}(\{A, B, C\})$. Choosing a canonical triangular decomposition for the radical \mathcal{I} with the variable ordering $X_1 \prec \dots \prec X_n$ is achieved by the following combinatorial construction. We refer to [7] for a more formal definition.

Definition 2.7. Consider a zero-dimensional variety V and denote by $\pi = \pi_{n-1}^n$ the projection which removes the last coordinate. To a point x in V , we associate $N(x) = \#\pi^{-1}(\pi(x))$, that is, the number of points lying in the same π -fiber as x . Then, we split V into the disjoint union $V_1 \cup \dots \cup V_d$, where for all $i = 1, \dots, d$, V_i equals $N^{-1}(i)$, that is, the set of points $x \in V$ which have $N(x) = i$. This splitting process is applied recursively to all varieties V_1, \dots, V_d , taking into account the fibers of the successive projections π_i^n , for $i = n-1, \dots, 1$. In the end, we obtain a family of pairwise disjoint, equiprojectable varieties, whose reunion equals V ; they form the *equiprojectable decomposition* of V .

3. Approximate Equiprojectable Decomposition in Dimension Zero

In this section, we consider a zero-dimensional variety $V \subseteq A^n(\mathbb{C})$ over \mathbb{Q} . Each point of V is given by approximate coordinates in a sense that we make precise in Definition 3.1. We aim at defining and computing an *approximate triangular decomposition* of V . To do so, we extend the construction given by Definition 2.7 and introduce a notion of an *approximate equiprojectable decomposition* of V in Definition 3.7. Then, to each approximate equiprojectable component, we associate an approximate triangular set, leading to Definition 3.8 of an *approximate triangular decomposition* of V .

Therefore, an approximate triangular decomposition of V is obtained by interpolating the points of V given by approximate coordinates. We provide stability analysis for this interpolation in Section 4. Moreover, we report on experiments that illustrate the accuracy of our stability analysis in Sections 5 and 6.

Definition 3.1. Let $\epsilon > 0$ and $r \geq 0$ be real numbers. Let $\bar{x} = (\bar{x}_1, \dots, \bar{x}_n)$ be a point of V and let $x = (x_1, \dots, x_n) \in A^n(\mathbb{C})$ with $x \neq 0$. We say that (x, r) is an

approximate point for \bar{x} with tolerance ϵ , denoted by $\bar{x} \simeq_\epsilon(x, r)$, if the following conditions hold for all $1 \leq i \leq n$:

- (i) $|\bar{x}_i - x_i| \leq r$,
- (ii) $r \leq \epsilon |x|$.

where $|x| = \max(|x_1|, \dots, |x_n|)$.

With the notations of Definition 3.1 let (x, r) be an approximate point for \bar{x} with tolerance ϵ . Let $1 \leq i \leq n$ be fixed. If \bar{x}_i and x_i are complex numbers and $\bar{x}_i \neq 0$ then a frequently-used measure of the number of correct significant decimal digits in the approximate coordinate x_i is the *logarithm of the relative error* $\text{lre}(x_i, \bar{x}_i)$ given by

$$\text{lre}(x_i, \bar{x}_i) = -\log_{10} \frac{|\bar{x}_i - x_i|}{|\bar{x}_i|}. \quad (3.1)$$

Properties (i) and (ii) of Definition 3.1 lead to

$$\text{lre}(x_i, \bar{x}_i) \geq -\log_{10} \epsilon - \log_{10} \frac{|x|}{|\bar{x}_i|}. \quad (3.2)$$

In practice, one requires $\epsilon < 1$ and thus Formula (3.2) gives a good measure of the approximation of coordinate \bar{x}_i by means of coordinate x_i . Similarly, Formula (3.3) below gives a good measure of the approximation of point \bar{x} by means of point x , for $x \neq 0$:

$$\text{lre}(x, \bar{x}) = -\log_{10} \frac{|\bar{x} - x|}{|\bar{x}|}. \quad (3.3)$$

As we shall see now, another good measure of this approximation is

$$\text{lb}(\bar{x}, x) = -\log_{10} \frac{|\bar{x} - x|}{|x|}. \quad (3.4)$$

Indeed, one can easily check that the following holds:

$$\left| \log_{10} \frac{|\bar{x} - x|}{|x|} - \log_{10} \frac{|\bar{x} - x|}{|\bar{x}|} \right| = \left| \log_{10} \frac{|\bar{x}|}{|x|} \right|. \quad (3.5)$$

Moreover, we claim that for all $\epsilon > 0$:

$$\left| \log_{10} \frac{|\bar{x}|}{|x|} \right| \approx \epsilon, \quad (3.6)$$

Thus, $\text{lre}(x, \bar{x})$ and $\text{lb}(\bar{x}, x)$ are very close when ϵ is very small. To prove our claim, we start from

$$||\bar{x}| - |x|| \leq |\bar{x} - x| \leq \epsilon |x|, \quad (3.7)$$

which holds by assumption (points (i) and (ii) of Definition 3.1). We deduce

$$\left| \frac{|\bar{x}|}{|x|} - 1 \right| \leq \epsilon. \quad (3.8)$$

Since ϵ is meant to be very small, using $\log_{10}(1 - \epsilon) \approx -\epsilon$ and $\log_{10}(1 + \epsilon) \approx \epsilon$, we finally obtain Formula (3.6).

A representation (using approximate points in the sense of Definition 3.1) of the isolated roots of the variety $V \subseteq A^n(\mathbb{C})$ of an input polynomial system $F = \{F_1, \dots, F_n\} \subset \mathbb{Q}[X_1, \dots, X_n]$ can be obtained by numerical homotopy construction. In particular, we used the PHC software [30]. Indeed, for each point \bar{x} of V , the corresponding solution x returned by PHC is given with the condition number of the Jacobian matrix of F at x , denoted by $cond$. The value $cond$ can be used to estimate the distance between \bar{x} and x (see [18] for details). More precisely, because we use double precision floating-point numbers in the computation, a reasonable formula is: $|\bar{x}_i - x_i| / |x_i| \approx cond \cdot 10^{-16}$ for all $1 \leq i \leq n$ (see Table 4). Given $\epsilon > 0$, with this estimate, one can check whether each isolated point \bar{x} of V admits approximate points within tolerance ϵ . Theoretically, the homotopy continuation method can obtain approximate points arbitrarily close to the exact roots for any tolerance ϵ . So, if the multiplicity of each point is 1, a one-to-one map between approximate roots and exact ones can be computed. Note that none of the systems used in Section 6 have multiple roots (see Table 2).

Remark 3.2. Let $\epsilon > 0$. From now on, we assume that for each point $\bar{x} \in V$ we are given $x \in A^n(\mathbb{C})$ and $r > 0$, such that $\bar{x} \simeq_\epsilon(x, r)$ holds. Then, we denote by \tilde{V} the set of all (x, r) , and we write $V \simeq_\epsilon \tilde{V}$.

We now return to the construction given by Definition 2.7. Again let $\pi = \pi_{n-1}^n$ be the natural projection from $A^n(\mathbb{C})$ to $A^{n-1}(\mathbb{C})$ which removes the last coordinate. Given two points \bar{x} and \bar{x}' of V we have to decide whether they lie in the same π -fiber. Since \bar{x} and \bar{x}' are given by approximate points we need the following.

Definition 3.3. Let i and j be integers such that $1 \leq i \leq j \leq n$. Let $\bar{x}, \bar{y} \in \pi_j^n(V)$. Let $x = (x_1, \dots, x_j)$ (resp. $y = (y_1, \dots, y_j)$) and (x, r) (resp. (y, r')) be approximate coordinates of \bar{x} (resp. \bar{y}) with tolerance ϵ . We say that \bar{x} and \bar{y} lie approximately in the same π_i^j -fiber with tolerance ϵ if for all $1 \leq k \leq i$ we have

$$|x_k - y_k| \leq r + r'. \quad (3.9)$$

Proposition 3.4. *With the notations of Definition 3.3, if the points $\bar{x}, \bar{y} \in \pi_j^n(V)$ are in the same π_i^j -fiber, that is, if $\pi_i^j(\bar{x}) = \pi_i^j(\bar{y})$ then, the points \bar{x} and \bar{y} lie approximately in the same π_i^j -fiber with tolerance ϵ .*

PROOF. By contradiction. Assume that \bar{x} and \bar{y} do not lie approximately in the same π_i^j -fiber with tolerance ϵ . Then, there exists $1 \leq k \leq i$ such that $|x_k - y_k| > r + r'$. By assumption, we have $\bar{x}_k = \bar{y}_k$, $|\bar{x}_k - x_k| < r$ and $|\bar{y}_k - y_k| < r'$. This leads to

$$|x_k - y_k| \leq |\bar{x}_k - x_k| + |\bar{y}_k - y_k| \leq r + r'. \quad (3.10)$$

A contradiction. \square

Remark 3.5. Suppose $1 \leq i \leq j \leq n$. For the points of $\pi_j^n(V)$, the relation “lying approximately in the same π_i^j -fiber with tolerance ϵ ” may not be an equivalence relation, because of roundoff errors. We need to exclude this situation in order to

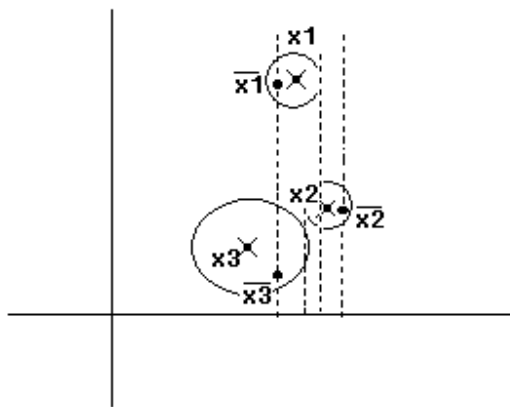


FIGURE 1. $\bar{x}_1, \bar{x}_2, \bar{x}_3$ are exact points, x_1, x_2, x_3 are the approximate points respectively. Here, \bar{x}_1, \bar{x}_2 lie in different fibers, but are approximately in the same fiber, and \tilde{V} satisfies the equivalence condition.

adapt the construction of Definition 2.7 with approximate points for the points of V . In theory, this situation may be avoided by reducing the tolerance ϵ , and thus the radius r at each point of V . However, in practice, for some systems it is hard to obtain the approximate roots when ϵ is very small. For example, for systems possessing a cluster of points for which we could not get a tolerance small enough, we would not be able to meet the requirements of Definition 3.7. These precautionary remarks being made, we will propose in Definition 3.7 a notion of an *approximate equiprojectable decomposition* of V , where the points of V are given by approximate points in the sense of Definition 3.1.

Definition 3.6. We say that \tilde{V} satisfies *the weak equivalence condition with tolerance ϵ* if for all $1 \leq i \leq j \leq n$, the relation "lying approximately in the same π_i^j -fiber with tolerance ϵ " is an equivalence relation in $\pi_j^n(V)$. Furthermore, we say that \tilde{V} satisfies *the strong equivalence condition with tolerance ϵ* if for every $\bar{x}, \bar{y} \in V$ with approximate points $(x, r), (y, r') \in \tilde{V}$, with tolerance ϵ , for all $1 \leq j \leq n$ the following conditions are equivalent:

- we have $\pi_j^n(\bar{x}) = \pi_j^n(\bar{y})$,
- the points \bar{x} and \bar{y} lie approximately in the same π_j^n -fiber.

Here we illustrate Definition 3.6 through Figures 1, 2 and 3 where we consider different \tilde{V} 's for the same V . In Figure 1, the set \tilde{V} satisfies the weak equivalence condition; observe that \bar{x}_1, \bar{x}_2 lie approximately in the same fiber, but \bar{x}_1 and \bar{x}_2 lie in different fibers. In Figure 2, the points \bar{x}_1, \bar{x}_2 and \bar{x}_1, \bar{x}_3 are pairs of points lying approximately in the same fiber, but \bar{x}_2, \bar{x}_3 do not lie approximately in the same

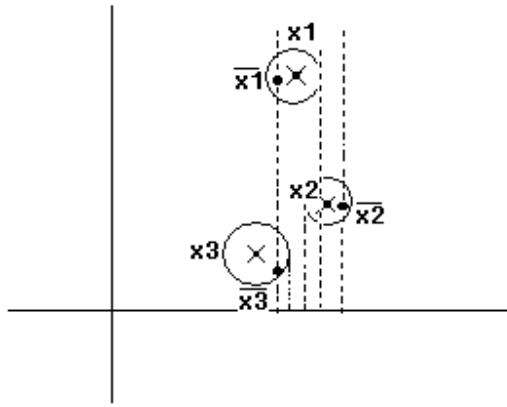


FIGURE 2. refining x_3 we get a smaller radius. Here, both pairs \bar{x}_1, \bar{x}_2 and \bar{x}_1, \bar{x}_3 lie approximately in the same fiber, but \bar{x}_2, \bar{x}_3 do not lie approximately in the same fiber. The set \tilde{V} does not satisfy the weak equivalence condition.

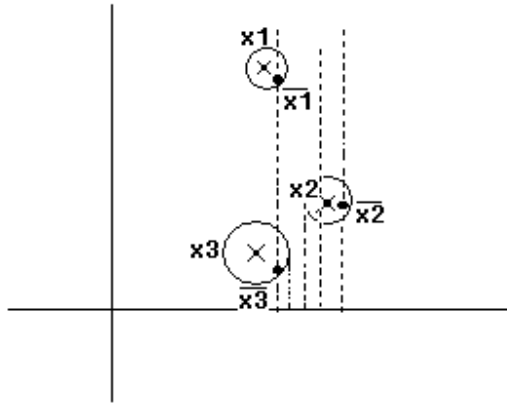


FIGURE 3. refining x_1 we get the correct result. Here, \bar{x}_1, \bar{x}_2 lie in different fibers and both weak and strong equivalence conditions are satisfied.

fiber. Hence, in this case, the set \tilde{V} does not satisfy weak equivalence condition. In Figure 3, we refine the three approximate roots until the weak equivalence condition is satisfied again (the strong equivalence condition is also satisfied); we see that \bar{x}_1, \bar{x}_2 lie in the different fibers.

In practice, the "exact" points of V are unknown, so we cannot determine whether the strong equivalence condition is satisfied or not. But we can detect whether the weak equivalence condition holds or not. However, in our experiments reported in Section 6, for each variety V , the exact points are known and we could decide whether \tilde{V} satisfies the strong equivalence condition.

If the weak equivalence condition is satisfied but the strong equivalence condition is not (e.g. see Figure 1), then there exists two distinct points $\bar{x}, \bar{y} \in V$, with respective approximate points $(x, r), (y, r')$, and an index $1 \leq i \leq n$ such that \bar{x}_i and \bar{y}_i are different but very close to each other; more precisely $|\bar{x}_i - \bar{y}_i| < 2r + 2r'$ holds (generally the distance $|\bar{x}_i - \bar{y}_i|$ will be less than 10^{-13} , see Table 4). Due to roundoff errors in numerical computation, we cannot always avoid these rare cases.

Finally, we note that introducing the notion of "weak equivalence condition" is needed by Definition 3.7.

Definition 3.7. Assume that \tilde{V} satisfies the weak equivalence condition with tolerance ϵ . Define $\pi = \pi_{n-1}^n$. To every point \bar{x} in V , we associate $N(\bar{x})$ the number of points in V which lie approximately in the same π -fiber as x with tolerance ϵ . For all $i \geq 1$, we denote by V_i the set of points $x \in V$ satisfying $N(x) = i$. Then, we split V into a disjoint union $V_1 \cup \dots \cup V_d$, for some $d \in \mathbb{N}$ large enough. This splitting process is applied recursively to all V_1, \dots, V_d , taking into account the fibers of the successive projections π_i^n , for $i = n-1, \dots, 1$. In the end, we obtain a family of pairwise disjoint subsets of V , whose union equals V ; they form an *approximate equiprojectable decomposition* of V with tolerance ϵ . If this approximate equiprojectable decomposition of V (with tolerance ϵ) consists of only one subset, that is, V itself, we say that V is *equiprojectable* with tolerance ϵ , otherwise the parts of the approximate equiprojectable decomposition of V (with tolerance ϵ) are called *approximate equiprojectable components* of V with tolerance ϵ .

Note that each approximate equiprojectable component of V is equiprojectable with tolerance ϵ . To each approximate equiprojectable component of V with tolerance ϵ we can associate an *approximate triangular set* by means of Definition 3.8. This leads to a notion of an *approximate triangular decomposition* for the variety V .

Definition 3.8. Assume that the zero-dimensional variety V is equiprojectable with tolerance ϵ . Then, by means of the interpolation formulas of Proposition 2.5 one can compute a triangular set $\{N_1, \dots, N_n\}$ called an *approximate triangular set* of V with tolerance ϵ .

Now, assume that V is not approximately equiprojectable with tolerance ϵ . A family of approximate triangular sets of approximate equiprojectable components of V (with tolerance ϵ) forms an *approximate triangular decomposition* of V , with tolerance ϵ .

4. Stability Analysis

In this section, we explore the relation between the *relative error* on the coordinates of the approximate points of V and the *relative error* on the interpolated polynomials of the approximate triangular decomposition given by Definition 3.8. The coefficients of a polynomial continuously depend on its roots. However, a small error in a root may result in a large error in the coefficients, motivating some of stability analysis.

For the relation between the errors mentioned above to be useful in practice, we must face the following fact: the relative error of a root cannot be computed when the exact root is unknown. In order to overcome this difficulty, for a point \bar{x} of V given by an approximate point (x, r) , we view the exact coordinates $\bar{x} = (\bar{x}_1, \dots, \bar{x}_n)$ as a random variable which takes values in the region defined by the following: for all $1 \leq i \leq n$

$$|x_i - \bar{x}_i| \leq r. \quad (4.1)$$

In this paper, we used the word *bias* instead of *relative error* in order to avoid conflicting terminology.

Definition 4.1. For $\bar{x}, x \in \mathbb{C}$, we call the the *bias* of x w.r.t. \bar{x} the fraction

$$\delta_x = \frac{\bar{x} - x}{x} \quad (4.2)$$

simply denoted by δ , when no confusion may occur.

Remark 4.2. We would like to observe at this point that none of the results of this section require knowledge of the exact coordinates of the points of V . Hence, our results apply also in practice to the situation where V is initially given by a polynomial system with inexact coefficients rather than a polynomial system with exact coefficients. Note that the PHC software [30, 18] can process both types of polynomial systems.

We define now the bias for the coefficients of a polynomial. Our definition applies to univariate polynomials as well as to multivariate polynomials. Let $e = (e_1, \dots, e_n) \in \mathbb{N}^n$ be an exponent vector. We denote by X^e the monomial $X_1^{e_1} \cdots X_n^{e_n}$ of $\mathbb{C}[X_1, \dots, X_n]$. We write $p = \sum_{e \in S} f_e X^e$ a polynomial of $\mathbb{C}[X_1, \dots, X_n]$ with (finite) support S . For every $e \in \mathbb{N}^n$ with $e \notin S$ we set to zero the coefficient f_e , i.e. we define $f_e = 0$. Hence we can simply write $p = \sum_e f_e X^e$.

Typically, in our stability analysis, the polynomial f of Definition 4.3 will be a polynomial interpolating the approximate coordinates of the points of V , whereas \bar{f} will be the corresponding polynomial obtained from the exact coordinates of the points of V .

Definition 4.3. Let $\bar{p} = \sum_e \bar{f}_e X^e$ and $p = \sum_e f_e X^e$ be polynomials in $\mathbb{C}[x_1, \dots, x_n]$. For every $e \in \mathbb{N}^n$, the *bias of coefficient* f_e w.r.t. \bar{p} is defined by

$$\delta_e = \frac{\bar{f}_e - f_e}{f_e}. \quad (4.3)$$

The *bias of the polynomial p* w.r.t. \bar{p} is the bias of the coefficient of p w.r.t. \bar{p} which has the largest norm.

The interpolated polynomials given by Proposition 2.5 are multivariate polynomials that are constructed as univariate polynomials over a suitable coefficient ring. Because of these formulas, we can focus on the univariate case. Let $\bar{p} \in \mathbb{C}[X]$ be a univariate monic polynomial of degree b given by approximate values x_1, \dots, x_b of its roots with respective radii r_1, \dots, r_b .

$$p = \prod_{i=1}^{i=b} (x - x_i). \quad (4.4)$$

Let $\delta_1, \dots, \delta_b$ be the respective bias of x_1, \dots, x_b such that the exact roots of \bar{p} are $x_1 + x_1 \delta_1, \dots, x_b + x_b \delta_b$. Hence we have

$$\bar{p} = \prod_{i=1}^{i=b} (x - x_i - x_i \delta_i). \quad (4.5)$$

Notation 1. *In the remainder of this section, we assume that $\delta_1, \dots, \delta_b$ are independent random (complex) variables, each of them with uniform distribution in a disk centered at 0 and with respective radii $r_1/|x_1|, \dots, r_b/|x_b|$. We define the bias bound and we denote it by ρ the maximum of $r_1/|x_1|, \dots, r_b/|x_b|$.*

In the proofs of Propositions 4.4, 4.6, and 4.7, we will denote by $O(\delta^2)$ any term in $\delta_i \delta_j$. When ρ is very small, we can ignore such higher order terms keep only the linear terms.

We will consider the bias of the polynomial \bar{p} w.r.t. p as a random variable denoted by γ . We direct the reader to the Appendix, for a brief review of the standard probability results we are using.

There are essentially three steps in computing the interpolated polynomials of Proposition 2.5:

- (I1) compute the univariate polynomials $e_{\alpha,i}$,
- (I2) compute the multivariate polynomials E_α , which are products of univariate polynomials $e_{\alpha,i}$,
- (I3) compute the multivariate polynomials N_ℓ which are sums of some multivariate polynomials.

For each step, we provide properties on the stability analysis of the corresponding calculations. For our study of the relation between \bar{p} and p , we need the following notation.

Notation 2. For $1 \leq k \leq b$, the k -th elementary symmetric function of x_1, \dots, x_b is given by

$$\sigma^k = \sum_{1 \leq a_1 < a_2 < \dots < a_k \leq b} x_{a_1} \cdots x_{a_k}, \quad (4.6)$$

and let $\sigma^0 := 1$. Observe that we have:

$$p = \prod_{i=1}^b (x - x_i) = \sum_{k=0}^b (-1)^k \sigma^k x^{b-k}. \quad (4.7)$$

Let $1 \leq j \leq b$. We denote by σ_j^k the element of $\mathbb{C}[x_1, \dots, x_n]$ obtained from σ^k by specializing x_j to 0, that is $\sigma_j^k = \sigma^k|_{x_j=0}$. Let l_j be the j -th Lagrange interpolation polynomial. Observe that we have:

$$l_j = \prod_{i=1, i \neq j}^b (x - x_i) = \sum_{k=0}^{b-1} (-1)^k \sigma_j^k x^{b-k-1}. \quad (4.8)$$

Proposition 4.4. *The bias γ of p w.r.t \bar{p} is bounded by*

$$\max \left(\frac{\sum_{i=1}^b |\sigma_i^k x_i|}{|\sigma^{k+1}|}, k = 0, \dots, b-1 \right) \rho. \quad (4.9)$$

We define

$$\varpi_k = \frac{\sqrt{3 \sum_{i=1}^b |\sigma_i^k x_i|^2}}{3 |\sigma^{k+1}|} \rho \quad (4.10)$$

$$\omega = \max(\varpi_k, k = 0, \dots, b-1). \quad (4.11)$$

If b is big enough, then γ is bounded by the normal distribution $N(0, \omega)$. (For the precise meaning of the statement being bounded by a distribution, please refer to Definition 7.3 in the Appendix.)

PROOF. By the definitions of \bar{p} and p , we have

$$\begin{aligned} \bar{p} - p &= \prod_{i=1}^b (x - x_i - x_i \delta_i) - \prod_{i=1}^b (x - x_i) \\ &= \prod_{i=1}^b (x - x_i) - \sum_{i=1}^b \prod_{j=1, j \neq i}^b (x - x_j) x_i \delta_i + O(\delta^2) - \prod_{i=1}^b (x - x_i) \\ &= - \sum_{i=1}^b l_i x_i \delta_i + O(\delta^2) \\ &\approx - \sum_{i=1}^b \left(\sum_{k=0}^{b-1} (-1)^k \sigma_i^k x_i \delta_i \right) x^{b-k-1} \\ &= - \sum_{k=0}^{b-1} (-1)^k \left(\sum_{i=1}^b \sigma_i^k x_i \delta_i \right) x^{b-k-1}, \end{aligned}$$

and

$$p = \prod_{i=1}^b (x - x_i) = \sum_{k=-1}^{b-1} (-1)^{k+1} \sigma^{k+1} x^{b-k-1}.$$

Thus, the absolute value of the bias for each coefficient γ_k , for $k = 0, \dots, b-1$, is given by

$$|\gamma_k| = \frac{|\sum_{i=1}^b \sigma_i^k x_i \delta_i|}{|\sigma^{k+1}|} \leq \frac{\sum_{i=1}^b |\sigma_i^k x_i|}{|\sigma^{k+1}|} \rho.$$

Hence, to order $O(\delta^2)$

$$\gamma \leq \max \left(\frac{\sum_{i=1}^b |\sigma_i^k x_i|}{|\sigma^{k+1}|}, k = 0, \dots, b-1 \right) \rho.$$

Recall that, by assumption, the random variables $\delta_1, \dots, \delta_b$ are independent. Also observe that, to order $O(\delta^2)$, the bias of each coefficient of p is a linear combination of these variables. Hence, we can compute the variance ω_k^2 of the bias γ_k of the coefficient x^{b-k-1} , for $k = 0, \dots, b-1$, by means of the properties given in the Appendix:

$$\begin{aligned} \omega_k^2 &= \text{Var} \left(\sum_{i=1}^b \sigma_i^k x_i \delta_i / \sigma^{k+1} \right) \\ &= \text{Var} \left(\sum_{i=1}^b \sigma_i^k x_i \delta_i \right) / |\sigma^{k+1}|^2 \\ &= \frac{\sum_{i=1}^b |\sigma_i^k x_i|^2}{|\sigma^{k+1}|^2} \text{Var}(\delta_i) \\ &\leq \frac{\sum_{i=1}^b |\sigma_i^k x_i|^2}{3|\sigma^{k+1}|^2} \rho^2 \\ &= \varpi_k^2. \end{aligned}$$

When b is big enough, the distribution of γ_k will tend to a normal distribution $N(0, \omega_k)$, by the results in the Appendix. Let $\omega = \max(\varpi_k, k = 0, \dots, b-1)$, then γ_k is bounded by $N(0, \omega)$ for each k . Finally, γ is bounded by $N(0, \omega)$. \square

Remark 4.5. If γ follows the normal distribution $N(0, \omega)$ and $x = 2\omega$ then we have $P(|\gamma| < x) \approx 0.95$. In fact, our experiments show that for $b \geq 10$, the probability $P(|\gamma| < x)$ is close to 0.95. Thus we can use Formula (4.10) to estimate the bias in the coefficients even if b is not very big. From the output of PHC we can estimate δ using condition numbers, compute ω , and finally estimate the bias for the coefficients with confidence level 0.95. In this section assuming b is big enough, then we have:

Proposition 4.6. Given n univariate polynomials, $p_i(x_i) = \sum_k a_{i,k} x_i^k$, $i = 1, \dots, n$, if each δ_i (the bias of p_i) satisfies $N(0, \omega)$, then the bias of $\prod_{i=1}^n p_i$ is bounded by $N(0, \sqrt{n}\omega)$ to order $O(\delta^2)$.

PROOF. Write the product of the univariate polynomials as a sum of monomials :

$$p_1 \cdots p_n = \sum f_e X^e,$$

where

$$f_e = f_{e_1, \dots, e_n} = a_{1, e_1} \cdots a_{n, e_n}.$$

Denote the exact coefficient by

$$\bar{f}_e = (a_{1, e_1} + a_{1, e_1} \delta_1) \cdots (a_{n, e_n} + a_{n, e_n} \delta_n).$$

By the same arguments as above:

$$\begin{aligned} \gamma_e &= \frac{\bar{f}_e - f_e}{f_e} \\ &= \frac{a_{1, e_1} \cdots a_{n, e_n} (\delta_1 + \cdots + \delta_n)}{a_{1, e_1} \cdots a_{n, e_n}} + O(\delta^2) \\ &\approx \delta_1 + \cdots + \delta_n. \end{aligned}$$

Because each δ_i satisfies $N(0, \omega)$, their sum is also normally distributed (see the Appendix) with distribution function $N(0, \sqrt{n}\omega)$. So, to order $O(\delta^2)$ the bias of $\prod_{i=1}^n p_i$ is bounded by $N(0, \sqrt{n}\omega)$. \square

Proposition 4.7. *Let $p_i(X) = \sum f_{i,e} X^e$, $i = 1, \dots, N$, be multi-variate polynomials such that δ_i (the bias of p_i) is normally distributed with distribution $N(0, \omega)$. Let*

$$\begin{aligned} \omega_e &= \frac{\sqrt{\sum_{i=1}^N f_{i,e}^2}}{|\sum_{i=1}^N f_{i,e}|} \omega \\ \omega' &= \max(\omega_e). \end{aligned} \tag{4.12}$$

Then, to order $O(\delta^2)$, the random variable γ for $\sum_{i=1}^N p_i(X)$ is bounded by $N(0, \omega')$.

PROOF. Examine the coefficients of the monomials:

$$\begin{aligned} p_1 + \cdots + p_N &= \sum f_e X^e \\ f_e &= f_{1,e} + \cdots + f_{N,e}. \end{aligned}$$

Let the exact coefficient be denoted by

$$\bar{f}_e = (f_{1,e} + f_{1,e} \delta_1) + \cdots + (f_{N,e} + f_{N,e} \delta_N).$$

Again, by the same arguments, the bias γ_e is:

$$\frac{\bar{f}_e - f_e}{f_e} = \frac{f_{1,e} \delta_1 + \cdots + f_{N,e} \delta_N}{f_{1,e} + \cdots + f_{N,e}} + O(\delta^2).$$

Because each δ_i is normally distributed by $N(0, \omega)$, the distribution of γ_e is still normal and equal to $N(0, \omega_e)$ (see the Appendix). So γ for the sum is bounded by $N(0, \omega')$ (again, see the Appendix for the meaning of bounded here). \square

Definition 4.8. Given an approximate triangular set T and the bias bound ρ of the approximate roots, let the *bias* of T be bounded by $N(0, \omega)$. Denote the *standard deviation* of T by sd where $sd = \omega/\rho$.

Remark 4.9. Let $V \simeq_{\epsilon} \tilde{V}$. Assume that \tilde{V} satisfies the strong equivalence condition with tolerance ϵ , in the sense of Definition 3.6. Then, it follows from Propositions 4.4, 4.6, and 4.7 that we can determine sd and the bias of the approximate triangular sets (in the approximate equiprojectable decomposition) of \tilde{V} with a given probability. Moreover, for an approximate system, given a perturbation of the approximate roots, we can estimate the change of the coefficients of the associated approximate triangular sets.

5. An illustrative example

Here we use a simple example to illustrate concept of approximate triangular set and our algorithm for determining the standard deviation. Let us consider:

$$sys = [zx^2 - zy, x^2 - 4y + y^2 + 2, -3zy + zy^2 + 3z - 3]. \quad (5.1)$$

The exact triangular set of the system with order $z \prec y \prec x$:

$$[z - 3, y^2 - 3y + 2, x^2 - y]. \quad (5.2)$$

1. Solving the system by PHC, we get 4 isolated points:

$$[z = 3.0, y = 2.0, x = 1.41421356237309, rco = 0.01511]$$

$$[z = 3.0, y = 1.0, x = 1.0, rco = 0.02089]$$

$$[z = 3.0, y = 2.0, x = -1.41421356237309, rco = 0.01511]$$

$$[z = 3.0, y = 1.0, x = -1.0, rco = 0.02089].$$

Here rco is the inverse of the condition number of Jacobian matrix at this point.

2. We remark, as we did in the Introduction, that each solved form $[z = 3.0, y = 2.0, x = 1.41421356237309]$, $[z = 3.0, y = 1.0, x = 1.0]$, $[z = 3.0, y = 2.0, x = -1.41421356237309]$, $[z = 3.0, y = 1.0, x = -1.0]$ is an approximate triangular set.
3. We use the condition numbers to estimate d_{\max} : $\delta = 1/rco \times 10^{-16} = 6.62 \times 10^{-15}$ and call this the estimated value of ρ . For this example, we know the exact solutions, and the exact distance between roots. In particular ρ should be $\sqrt{2} - 1.41421356237309 = 5.1 \times 10^{-15}$. In practice we don't know the exact solution of input system, and we only can give an estimated value for ρ . But we need to point out that this estimation works well for many examples. Comparisons are given in next section.
4. By the definition of approximate equiprojectable decomposition, the projection of the first and third points above are numerically equal since $|2.0 - 2.0| < (2.0/0.01511 + 2.0/0.01511) \times 10^{-16}$.

Also the projection of first and second points are not numerically equal since $|2.0 - 1.0| > (2.0/0.01511 + 1.0/0.02089) \times 10^{-16} = 1.8 \times 10^{-14}$.

In the same way, we get two different projected points $p1 = (3.0, 2.0)$, $p2 = (3.0, 1.0)$ on zy -plane, and there are two points on each fiber. The projections of $p1$, $p2$ onto the z axis is just one point $z = 3.0$. So the variety of sys is approximately equiprojectable. From the cardinality of the fibers,

# roots	# tests	% of trials: rel. err. > 1 sd (0.32 expected)	% of trials: rel. err. > 2 sd (0.05 expected)	% of trials: rel. err. > 3 sd (0.003 expected)
10	1000	0.328	0.0503	0.0168
20	1000	0.312	0.0425	0.0050
30	1000	0.350	0.0579	0.0023
40	800	0.335	0.0517	0.0067
50	500	0.342	0.0474	0.0042

TABLE 1. Experiments for our probabilistic analysis (sd = standard dev., rel. err. is relative error)

we know the degree sequence is $[1, 2, 2]$ with respect to the main variables of each polynomial in the triangular set. The degree sequence can be equivalently written as $1 \cdot 2^2$.

5. By formula 2.7, we get the approximate triangular set of sys :

$$[-.999999999999986y + 1.0x^2, y^2 - 3.0y + 2.0, z - 3.0]. \quad (5.3)$$

The biggest relative error of coefficients is 1.4×10^{-14} . By formula (4.10) and (4.12) the standard deviation (sd) is 2.89.

So $sd \times \rho = 1.9 \times 10^{-14} > 1.4 \times 10^{-14}$ is a good estimate for the relative error. In the next section we will give more nontrivial examples to support our statement. Because of the input error and round off error in numerical computation, there will be some monomials of approximate triangular sets with very small coefficients that do not appear in exact triangular set. Then the biggest relative error of coefficients is 1. So in practice we will consider coefficients which are smaller than a given tolerance as 0.

6. Experimental Results

We have conducted two sets of experiments. The first one illustrates the probabilistic analysis of Proposition 4.4. Experiments are described in Section 6.1, and the results appear in Table 1.

The second set of experiments deal with the computation of exact and approximate triangular decompositions. Section 6.2 presents the exact case whereas Section 6.3 reports on the approximate one. Most of the test polynomial systems that we use (see Table 2) are well known problems [1, 7, 29]. They are zero-dimensional square systems defined by multivariate polynomials over \mathbb{Q} generating radical ideals. Table 3 shows data for the exact triangular decompositions of these systems, the output by PHC is collected in Table 4, and Table 5 shows their approximate triangular decompositions computed from the PHC output. The main results for the purposes of this paper are given by this latter table.

Sys	Name	n	d	h	H	\hat{H}	reference
1	Issac97	4	2	2	71	1498	[29]
2	L3	3	3	1	1	1678	[1]
3	Sendra	2	7	7	59	2421	[29]
4	fabfaux	3	3	13	72	2650	[10]
5	L4	3	4	1	2	3977	[1]
6	Cylohexne	3	4	3	9	4361	[29]
7	Weispfenning94	3	5	0	10	7392	[29]
8	UteshevBikker	4	3	3	88	7908	[29]
9	Fee-1	4	2	2	34	23967	[29]
10	Reimer-4	4	5	1	14	56013	[29]
11	S9 ₁	8	2	2	33	58116	[29]
12	eco6	6	3	0	12	105718	[29]
13	Geneig	6	3	2	82	114466	[29]
14	gametwo5	5	4	8	674	158075	[29]
15	dessin-2	10	2	7	436	360596	[29]
16	eco7	7	3	0	26	387754	[29]
17	Methan61	10	2	16	227	452756	[29]

TABLE 2. Input systems ($n = \#$ polys.; $d =$ degree system; $h =$ height input coeffs; $H =$ height output coeffs; $\hat{H} =$ estimated height output coeffs.)

Sys	Exact equiproj dec. tim. (secs)	Degree configuration	#C-roots	Time to isolate \mathbb{R} -roots (secs)	# \mathbb{R} -roots
1	164	$16\ 1^3$	16	< 1	0
2	< 1	(1 3 1), (8 1 1), (8 2 1)	27	< 1	5
3	33	$46\ 1$	46	5	6
4	28	$27\ 1^2$	27	1	3
5	1	(24 2 1), (16 1 1)	64	< 1	8
6	6	(4 1 2), (8 1 1)	16	< 1	12
7	72	$54\ 1^2\ 1$	54	< 1	0
8	29201	$36\ 1^3$	36	7	10
9	24	$26\ 1^3$	26	2	6
10	10097	$18\ 2\ 1^2$	36	5	4
11	26	$10\ 1^7$	10	1	4
12	50	$16\ 1^5$	16	< 1	4
13	18	$10\ 1^3$	10	2	10
14	24320	$44\ 1^4$	44	45	12
15	527	$1\ 42\ 1^8$	42	15	1
16	2742	$32\ 1^6$	32	4	8
17	6251	$27\ 1^9$	27	28	13

TABLE 3. Exact equiprojectable triangular decomposition with the RegularChains library

Sys	#C-roots	#C-roots by PHC	PHC tim.(secs)	estimated ρ	exact ρ
1	16	16	1	0.448e-14	0.239e-14
2	27	27	1	0.186e-14	0.337e-14
3	46	46	4	0.159e-11	0.274e-14
4	27	27	2	0.224e-14	0.154e-14
5	64	64	1	0.143e-14	0.331e-14
6	16	16	< 1	0.835e-14	0.181e-14
7	54	49	5	0.183e-13	0.336e-14
8	36	36	6	0.767e-12	0.781e-14
9	26	26	5	0.229e-11	0.759e-14
10	36	36	3	0.739e-13	0.544e-14
11	10	10	3	0.107e-13	0.125e-14
12	16	16	3	0.292e-13	0.287e-14
13	10	10	2	0.629e-13	0.105e-13
14	44	43	6	0.665e-12	0.144e-13
15	42	41	11	0.585e-7	0.271e-14
16	32	32	14	0.760e-13	0.264e-14
17	27	13	10	0.846e-6	0.563e-13

TABLE 4. Approximate roots by PHC where the *estimate* ρ = condition number $\times 10^{-16}$ and *exact* ρ = largest 2-norm of distance between exact and approx root divided by the 2-norm of approx root.

6.1. Normal distribution test

Let b be a number of roots given in the column *# roots*. We randomly generate b roots, and view them as the exact roots of a polynomial \bar{p} of degree d . Then, we perturb each of these roots by a uniformly distributed random variable, leading to an approximate polynomial p . The two polynomials \bar{p} and p are expanded in order to obtain ε , the largest relative error for a coefficient. We compute the standard deviation sd by formula (4.10), and compare it with ε . These experiments are repeated many times (between 500 and 1000, see the column *# tests*) for $b = 10, 20, 30, 40, 50$. The third column is the percentage of times for which the relative error is bigger than one standard deviation. If the relative error is normally distributed, then this percentage should be 0.32, which we verify in our tests.

6.2. Exact triangular decomposition

The test polynomial systems are given in Table 2. For each input system F , we give n the number of variables, d the total degree of F , the logarithm h of the largest coefficient, the number of digits H appearing in the largest coefficient in the (exact) equiprojectable decomposition of F , and the height \hat{H} of that coefficient as estimated by the formulas of [7].

In order to compute the exact equiprojectable decomposition, we use the `RegularChains` library written in MAPLE by Lemaire, Moreno Maza and Xie [16] in which the algorithms of [19, 7] are implemented. Our computations are done on a

Sys	sd	exact $\rho \cdot sd$	δ_{coeff}	$< sd?$	$< 2sd?$	residual
1	403.3	0.9639e-12	0.197e-12	yes	yes	0.444e-15
2	7.492	0.2529e-13	0.211e-13	yes	yes	0.125e-13
3	1729.2	0.4736e-11	0.542e-11	no	yes	0.89e-11
4	1056.7	0.1625e-11	0.463e-12	yes	yes	0.201
5	59188.4	0.1959e-09	0.248e-09	no	yes	0.555e-7
6	23835.5	0.4314e-10	0.179e-11	yes	yes	0.7e-13
7	NA	NA	NA	NA	NA	NA
8	383.8	0.2996e-11	0.942e-12	yes	yes	0.163e-8
9	151.6	0.1151e-11	0.181e-12	yes	yes	0.504e-13
10	3928.4	0.2137e-10	0.397e-12	yes	yes	0.193e-18
11	45.77	0.5708e-13	0.133e-13	yes	yes	0.188e-15
12	121.7	0.3488e-12	0.184e-12	yes	yes	0.216
13	551.7	0.5815e-11	0.761e-13	yes	yes	0.314e-17
14	NA	NA	NA	NA	NA	NA
15	NA	NA	NA	NA	NA	NA
16	317.7	0.8397e-12	0.154e-11	no	yes	0.218e20
17	NA	NA	NA	NA	NA	NA

TABLE 5. Approximate Triangular Sets: sd = standard dev. defined in Section 4; $exact \rho$ = largest 2-norm distance between exact and approx root divided by the 2-norm of approx root; δ_{coeff} = largest relative error of coeffs of approx triangular set compared with the exact one.

2799 MHz Pentium 4 machine. The timings of exact equiprojectable decomposition are given in the first column of Table 3. To understand these timings, we should mention that the `RegularChains` code is high-level interpreted code (and not compiled). Moreover, this code is not supported by fast arithmetic, such as FFT-based arithmetic.

Each degree configuration specifies the degree sequences of the triangular sets in the decomposition (see [1] for similar data). Hence, the number of sequences in a degree configuration equals the number of equiprojectable components of the system. In Table 3 $\#\mathbb{C}$ -roots and $\#\mathbb{R}$ -roots are respectively the total number of complex and real roots of the system. The column labeled Time to isolate \mathbb{R} -roots, gives the total time in seconds to isolate all the real roots to a precision of 2^{-30} using interval arithmetic.

We have also isolated each complex root. This was done by Éric Schost (École Polytechnique, France) using `Magma` as follows. First, the *splitting circle* method of Schönhage was used to separate the complex roots. Then, Newton iteration was used to refine the isolation boxes. A precision of 200 digits could be achieved for our 17 test systems in less than 10 minutes on a Pentium P3 running at 1Ghz.

6.3. Approximate triangular sets

We use the PHC package [30, 18] to compute the approximate isolated roots of all the benchmark systems. Then we interpolate the approximate triangular sets and give the results of our error analysis for each system. The computations in Tables 4 and 5 are done on a 1.5 GHz Pentium M machine, and the timings for finding the roots using PHC are listed in *PHC Timing* of Table 4. In Table 4: the first column is the exact number of roots and second column is the number of roots found by PHC. For some systems, PHC (in black box mode) does not get every root. This simply means that the default settings in the black box version of PHC did not solve the system. We did not compute the approximate triangular sets for such systems. Some of these systems could certainly have been solved by using PHCPack, by exploiting the flexibility of its powerful user specified options, designed for more challenging problems. But we did not do that here. The *estimate* ρ is defined as the condition number $\times 10^{-16}$, and *exact* $\rho = \max(|x_i - \bar{x}_i|/|x_i|)$, $\bar{x}_i \in V$ where the \bar{x}_i are the "exact" roots, the x_i are the roots given by PHC, and the distance is given by the 2 norm. The results show that our estimated distance is often larger than the exact distance.

In Table 5: The second column gives the standard deviation of the approximate triangular set, as discussed in Remark 4.9. The third column is the product of the exact ρ and one standard deviation. In the fourth column δ_{coeff} is the largest relative error of the coefficients of the approximate triangular set as compared with the exact one. If this relative error is less than exact $\rho \cdot sd$, the element of the fifth column (labeled $< sd?$) is "yes", otherwise it is "no". Moreover, for every approximate triangular set, the relative error is bounded by $2sd$ (see column 6). The last column, labeled *residual*, gives the maximum residual of an approximate triangular set at the roots given by PHC. The results of this table support the conclusions of Remark 4.9.

7. Discussion

Exact triangular decomposition of exact polynomial systems have proved valuable in applications. There are well-developed algorithms for computing them and considerable recent improvements in their time complexity [7]. Such representations are desirable, not only because of their triangular solved-form structure, but also because, in comparison to other exact methods, their space complexity is well controlled [8]. In particular, they use the minimum number of polynomials needed to describe the equi-dimensional decomposition components of a polynomial system.

In this paper we have extended such methods to approximate systems of polynomials in the dimension zero case.

We have exploited the methods of Sommese, Verschelde, and Wampler [25, 30, 26] and the newly developing area of Numerical Algebraic Geometry, together with new ideas of Dahan, Schost, Moreno Maza, Wu, and Xie for forming the so-called equiprojectable decomposition [6] of a zero-dimensional variety.

Throughout this paper we have assumed that the input is a zero dimensional radical ideal. We briefly discuss the situation where both restrictions are removed, that is we consider general input polynomial systems. The methods of Sommese, Verschelde, and Wampler yield isolated points, possibly of higher multiplicity corresponding to the 0 dimensional equi-dimensional components. Such multiplicities can be removed (deflated) numerically using the techniques of [9] and [17] (see [15] for a symbolic method for the exact case). Hence the methods of this paper can be applied.

Our contribution in the zero-dimensional case, has been to show that the isolated points, given by approximate coordinates, can be interpolated in order to obtain the triangular decomposition for the zero-dimensional components which is an approximation of the exact equiprojectable decomposition. The methods of Sommese, Verschelde, and Wampler yield a numerical irreducible decomposition for this case, in particular they give a collection of triangular sets, each of them corresponding trivially to an isolated point.

In addition the $n-1$ dimensional components can be numerically interpolated by the methods of Sommese, Verschelde, and Wampler to obtain single polynomial which can be considered as a representation for the highest dimensional components with triangular shape. In addition their methods also give (non-triangular) representations of all of the positive dimensional components using generic points on each component. The above results, together with those in our paper on linearized triangular decompositions [21], represent progress on the general problem of obtaining approximate triangular representations for all components of a given polynomial system. Our detailed study of the zero dimensional case is particularly important as a preparation for the study of the general case.

The interpolation formulas of Dahan and Schost are an extension of classical Lagrange interpolation formulas for univariate polynomials to multivariate polynomials. For improperly distributed interpolation points, such as uniformly spread points, the interpolation problem can be ill-conditioned regardless of whether one uses a Lagrange or a Newton formulation [3, 12]. One way to avoid unstable interpolation is to choose some specially distributed interpolation points, such as Chebyshev or Legendre points. Obviously this cannot apply in the context of this paper: indeed, the location of the interpolation points is fixed; we have no freedom to change them.

In [4], the authors compute an exact absolute factorization of a bivariate polynomial from an approximate factorization. It is natural to ask if one could compute an exact equiprojectable decomposition from an approximate one. One preliminary answer is as follows. Let F be an (exact) input polynomial system in $\mathbb{Q}[X_1 \prec \dots \prec X_n]$ with total degree d and largest coefficient h . Then, the height of any coefficient of any (exact) triangular set in the equiprojectable decomposition of $V(F) \subseteq A^n(\mathbb{C})$ is in $O(h n d^n)$ [7]. This suggests that the numbers d and n must be small for this *reconstruction* (from approximate to exact) to be realistic. However, the question remains open for future work. Indeed, Table 6 shows that

the actual coefficient size H in the triangular set is much less than the above height upper bound \hat{H} .

Standard deviation is a measure of error distribution of coefficients. The very large standard deviation means the coefficients are very sensitive to the change of roots. For such system, interpolation is not a good method to get the approximate triangular set from the roots. For example, recall system 5 in Table 6 of Section 6. Although the error of roots is $0.331 \cdot 10^{-14}$, the error of coefficients accumulates up to $0.248 \cdot 10^{-9}$ because of its large sd .

In [27] fundamental theorem on backward error are given for polynomials. We will use these results to exploit the backward error for approximate triangular sets in a future paper.

When the input system F is approximate, although discontinuous phenomena can occur, some continuity aspects remain under the perturbation [28]. The favorable properties of the equiprojectable decomposition of $V(F)$ under specialization [7] suggests that the continuity of approximate equiprojectable decomposition needs to be discussed in future work.

References

- [1] P. Aubry and M. Moreno Maza. Triangular sets for solving polynomial systems: A comparative implementation of four methods. *J. Symb. Comp.*, 28(1-2):125–154, 1999.
- [2] P. Aubry and A. Valibouze. Using Galois ideals for computing relative resolvents. *J. Symb. Comp.*, 30(6):635–651, 2000.
- [3] J.P. Berrut and L.N. Trefethen. Barycentric lagrange interpolation. *SIAM Rev.*, 2005. to appear.
- [4] G. Chèze and A. Galligo. From an approximate to an exact absolute polynomial factorization. Technical report, Université de Nice, 2005.
- [5] S.C. Chou. *Mechanical Geometry Theorem Proving*. D. Reidel Publ. Comp., Dordrecht, 1988.
- [6] X. Dahan, M. Moreno Maza, É. Schost, W. Wu, and Y. Xie. Equiprojectable decompositions of zero-dimensional varieties. In *ICPSS*, pages 69–71. University of Paris 6, France, 2004.
- [7] X. Dahan, M. Moreno Maza, É. Schost, W. Wu, and Y. Xie. Lifting techniques for triangular decompositions. In *ISSAC'05*, pages 108–115, ACM Press 2005.
- [8] X. Dahan and É. Schost. Sharp estimates for triangular sets. In *ISSAC 04*, pages 103–110. ACM Press, 2004.
- [9] Barry H. Dayton and Zhonggang Zeng. Computing the Multiplicity Structure in Solving Polynomial Systems. In *ISSAC 05*, pages 116–123. ACM Press, 2005.
- [10] European Commission. *FRISCO - A Framework for Integrated Symbolic/Numeric Computation*. Esprit Scheme Project No. 21 024, 1996. <http://www.nag.co.uk/projects/FRISCO.html>.

- [11] X.S. Gao and Y. Luo. A characteristic set algorithm for difference polynomial systems. pages 28–30, 2004.
- [12] N.J. Higham. The numerical stability of barycentric lagrange interpolation. *IMA Journal of Numerical Analysis*, 24(547-546), 2004.
- [13] M. Kalkbrener. A generalized euclidean algorithm for computing triangular representations of algebraic varieties. *J. Symb. Comp.*, 15:143–167, 1993.
- [14] D. Lazard. Solving zero-dimensional algebraic systems. *J. Symb. Comp.*, 13:117–133, 1992.
- [15] G. Lecerf. Quadratic Newton iteration for systems with multiplicity. *Found. Comput. Math.*, 2:247–293, 2002.
- [16] F. Lemaire, M. Moreno Maza, and Y. Xie. The RegularChains library. In *Maple 10*, Maplesoft, Canada. To appear.
- [17] Anton Leykin, Jan Verschelde, and Ailing Zhao. Evaluation of Jacobian Matrices for Newton’s Method with Deflation to approximate Isolated Singular Solutions of Polynomial Systems. In *SNC 2005 Proceedings*, pages 19–28.
- [18] A. Leykin and J. Verschelde. Phcmaple: A maple interface to the numerical homotopy algorithms in phcpack. In *proceedings of ACA’04*, pages 139–147, University of Texas at Beaumont, USA, 2004.
- [19] M. Moreno Maza. On triangular decompositions of algebraic varieties. Technical Report 4/99, NAG, UK, Presented at the MEGA-2000 Conference, Bath, UK, 2000. <http://www.csd.uwo.ca/~moreno>.
- [20] M. Moreno Maza and G. Reid and R. Scott and W. Wu. On Approximate Triangular Decompositions I. Dimension Zero. In D. M. Wang and L. Zhi editors. *Symbolic-Numeric Computation*, page 250-275, Xi’an, China, 2005.
- [21] M. Moreno Maza and G. Reid and R. Scott and W. Wu. On Approximate Triangular Decompositions II. Linear Systems. In D. M. Wang and L. Zhi editors. *Symbolic-Numeric Computation*, page 276-296, Xi’an, China, 2005.
- [22] J. F. Ritt. *Differential Equations from an Algebraic Standpoint*, volume 14. American Mathematical Society, New York, 1932.
- [23] É. Schost. Complexity results for triangular sets. *J. Symb. Comp.*, 36(3-4):555–594, 2003.
- [24] A. N. Shiryaev. *Probability*. Springer, 1995.
- [25] A.J. Sommese and J. Verschelde. Numerical homotopies to compute generic points on positive dimensional algebraic sets. *J. Complexity*, 16(3):572–602, 2000.
- [26] A.J. Sommese, J. Verschelde, and C.W. Wampler. Numerical decomposition of the solution sets of polynomial systems into irreducible components. *SIAM J. Numer. Anal.*, 38(6):2022–2046, 2001.
- [27] Hans J. Stetter. The nearest polynomial with a given zero, and similar problems. *SIGSAM Bull.*, 33(4):2–4, 1999.
- [28] Hans J. Stetter and Gunther H. Thallinger. Singular systems of polynomials. In *ISSAC ’98: Proceedings of the 1998 international symposium on Symbolic and algebraic computation*, pages 9–16, New York, NY, USA, 1998. ACM

- Press.
- [29] The symbolicdata project, 2000–2002. <http://www.SymbolicData.org>.
- [30] J. Verschelde. PHCpack: A general-purpose solver for polynomial systems by homotopy continuation. *ACM Transactions on Mathematical Software*, 25(2):251–276, 1999.
- [31] D. M. Wang. *Elimination Methods*. Springer, Wein, New York, 2000.
- [32] W. T. Wu. A zero structure theorem for polynomial equations solving. *MM Research Preprints*, 1:2–12, 1987.

Acknowledgment

The authors would like to thank Yuzhen Xie (University of Western Ontario, Canada) who realized all the experiments with the `RegularChains` library [16]. And we also appreciate the following colleagues for their helps: François Lemaire (University of Lille 1, France) who provided the source code of `Realroots` to compute the real solutions with `RegularChains` library and Éric Schost (École Polytechnique, France) who isolated each complex root by `Magma`.

Appendix - Brief review of probability theory

In our stability analysis of coefficients, a probability model was introduced. We will give a brief review of the relevant standard probability knowledge required.

- If δ is a random variable and c is a constant in \mathbb{R} then $Var(c\delta) = c^2Var(\delta)$.
- If $\delta_1, \dots, \delta_b$ are random variables and $\xi = \sum \delta_i$ then expectation value is additive: $E(\xi) = \sum E(\delta_i)$. Moreover if they are independent, then the variance of sum of these random variables is also additive: $Var(\xi) = \sum Var(\delta_i)$.
- Suppose $\delta = \delta_{re} + \delta_{im}\sqrt{-1}$ and δ_{re}, δ_{im} are independent random variables with the same distribution with $c \in \mathbb{C}$. Then $Var(\Re(c\delta)) = |c|^2Var(\delta_{re}) = Var(\Im(c\delta)) = |c|^2Var(\delta_{im})$, where $\Re(z)$ and $\Im(z)$ are the real and imaginary parts of z . In this paper we define $Var(\delta) := Var(\delta_{re})$.
- $N(0, 1)$ is standard normal distribution with mean 0, standard deviation 1, probability density function $p(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$ and cumulative density function $\Phi(x) = \int_{-\infty}^x p(x)dx$. Note that $\Phi(1) \approx 0.68$, $\Phi(2) \approx 0.95$.
- Suppose that $\delta_1, \dots, \delta_b$ are independent random variables with distribution functions F_1, \dots, F_b and $E(\delta_i) = 0$, $0 < Var(\delta_i) < \infty$, $s_b^2 = \sum Var(\delta_i)$. The Lindeberg condition for a sum of independent random variables is that for any $t > 0$:

$$\frac{1}{s_b^2} \sum_{k=1}^b \int_{|x| > ts_b} x^2 dF_k(x) \longrightarrow 0 \quad \text{when } b \longrightarrow \infty \quad (7.1)$$

From our assumptions about the roots, the bias is uniformly distributed and because $0 < Var(\delta_i) < \infty$ we have $s_b^2 \rightarrow \infty$ as $b \rightarrow \infty$. So for any $t > 0$, there always exists L , when $b > L$ the integral above is 0.

Proposition 7.1 (uniform distribution and Lindeberg condition). *If $\delta_1, \dots, \delta_b$ are independent random variables with uniform distribution, and $E(\delta_i) = 0$, if the variance of each δ_i is nonzero and finite then this family of random variables satisfies the Lindeberg condition.*

Proposition 7.2 (Lindeberg's central limit theorem [24]). *Suppose $\delta_1, \dots, \delta_b$ are uniformly distributed independent random variables, $E(\delta_i) = 0$ and δ_i satisfies the Lindeberg condition. Let $S_b = \sum_{i=1}^b \delta_i$ then when $b \rightarrow \infty$, the sum of variables divided by its standard deviation is convergent (in distribution) to a standard normal distribution:*

$$\frac{S_b}{s_b} \rightarrow N(0, 1) \quad \text{as } b \rightarrow \infty \quad (7.2)$$

Definition 7.3. We say a random variable ξ or $|\xi|$ is bounded by $N(0, \omega)$ if the probability $P(|\xi| < x\omega) > \Phi(x)$.

When ω is bigger, the probability will also be bigger. In particular if $\omega' > \omega$ then $P(|\xi| < x\omega') > P(|\xi| < x\omega)$, so ξ is also bounded by $N(0, \omega')$.

Marc Moreno Maza
e-mail: moreno@orcca.on.ca

Greg Reid
e-mail: reid@uwo.ca

Robin Scott
e-mail: rscott2@uwo.ca

Wenyuan Wu
e-mail: wwu25@uwo.ca

ORCCA
Computer Science Department
University of Western Ontario
London
Canada