

Analysis of Zero Clusters in Multivariate Polynomial Systems

Hans J. Stetter, Vienna

February 20, 1996

Abstract

We consider a cluster of m zeros of a multivariate polynomial system which we interpret as a perturbation of a system with an m -fold zero. By algebraic techniques, we find a first order correct representation of the primary ideal of the cluster zeros from which we obtain approximations for the individual zeros in the cluster.

1 Conceptual Background

Our considerations will proceed in the conceptual context of *numerical nonlinear algebra*: We will consider sets P of polynomials $p_\nu \in \mathbb{C}[x_1, \dots, x_s] =: \mathbb{P}^s$ whose coefficients have specified *numerical* values. Throughout, we will assume that $\langle P \rangle$ is a *0-dimensional ideal* and that we have a *complete intersection* case, with P containing exactly s polynomials. Our analysis is a contribution to the following overall task: *Find sufficiently good approximations for all zeros of P* . Note that some or all zeros of P will generally be *irrational* complex s -tuples, even if P has integer coefficients, so that a numerical specification of a zero will necessarily have to be an *approximation*. However, the concepts of *zero* and *sufficiently good* in the above sentence need further clarification. In agreement with standard practice in numerical mathematics we distinguish two cases:

Case 1: All coefficients in P may be assumed to be exact as specified. Here, we have the classical algebraic meaning of a zero:

$$z \in \mathbb{C}^s \text{ is a zero of } P \iff p_\nu(z) = 0, \nu = 1(1)s,$$

and an approximation z is *sufficiently good* if, for some appropriate norm in \mathbb{C}^s ,

$$\|(p_\nu(z))\| \leq \alpha \quad \text{for a specified } \alpha > 0.$$

Case 2: Some of the (non-vanishing) coefficients in P are only known to a specified level of accuracy. Then P represents an *equivalence class* \overline{P} of polynomial sets \tilde{P} ; the members of \overline{P} cannot be distinguished in the given context. Thus, the concept of a zero has to be widened to that of a *pseudozero*:

$$z \in \mathbb{C}^s \text{ is a pseudozero of } P \iff \exists \tilde{P} \in \overline{P}: p_\nu(z) = 0, p_\nu \in \tilde{P},$$

and z is a *sufficiently good approximation* if it is a pseudozero of P . Case 2 occurs predominantly in the simulation of real-world phenomena by mathematical models whose data are generally not known to an arbitrary accuracy.

Contrary to a first impression, this distinction – which is omnipresent in numerical computations – does *not* affect the design of solution procedures, due to a general Design Principle of Numerical Mathematics:

Step 1: Find a reasonable approximation of the solution efficiently and robustly.

Step 2: Verify the quality of the approximation; if not sufficient goto Step 3.

Step 3: Improve the present approximation; goto Step 2.

The distinction of the above cases 1 and 2 only appears in Step 2 whereas Steps 1 and 3 are independent of this distinction. Naturally, the potentially iterative improvement in Steps 2/3 needs an emergency exit: solution procedures may fail to give a satisfactory answer.

For the determination of a well-isolated zero z of P , the classical procedure in Step 3 is Newton's method. In the case of a sufficiently dense *cluster* of zeros, however, where the Jacobian of P is nearly rank-deficient over a large part of the region occupied by the cluster, this approach meets with great difficulties; without very accurate starting approximations, no or only a few of the zeros of the cluster may be found by Newton's method or some variant of it. Therefore, the main part of this paper will present an alternate approach for a localization of all individual cluster zeros.

Such an approach must take into account the fact that the task of localizing the zeros of a cluster is severely *ill-conditioned*: The accurate locations of the zeros are extremely sensitive to small changes in the data (i.e. the polynomial coefficients). However, the *location of the cluster* and its *multiplicity* (the number of zeros in the cluster) generally do not share this sensitivity; this has been well-known for some time, an early analysis has been given, e.g., by [1]. Actually, the occurrence of a cluster of m zeros (an m -cluster) is equivalent to the presence of an m -fold zero in a system whose data are close to those of the specified problem; the location of such an m -fold zero depends moderately on the specified data. Therefore, the appropriate version of Step 1 with respect to an m -cluster of zeros is:

Find the approximate location z_0 of an m -fold zero of a reasonably close system P_0 , by a stable and efficient floating-point computation.

While the realization of this task, e.g. on the basis of a stabilized floating-point Groebner basis computation, is in itself extremely interesting and presently under investigation, we will assume the result of this Step 1 as given for the purpose of this paper. Step 2 then amounts to the verification of the distance between P_0 and P . If the criteria specified for the particular situation are met, we are finished: Within the specified accuracy level, our system P cannot be distinguished from the system P_0 with a genuine m -fold zero at z_0 .

On the other hand, if P_0 is not sufficiently close to P by our standards, we must design a Step 3 procedure which *splits* the m -fold zero z_0 into m individual zeros and finds their approximate locations. This is the task to which we will now devote ourselves.

Since we will use floating-point arithmetic, the condition of this task depends essentially on the *relative* separation of the individual zeros in the cluster. This can be significantly improved by moving the “cluster center” z_0 to the origin. However, the computation of some of the coefficients in the shifted polynomials $\vec{p}_\nu(\xi) := p_\nu(z_0 + \xi)$ will meet with extreme cancellation of leading digits as these coefficients vanish in the shifted versions of the $p_{0\nu}$, with their m -fold zero at the origin. Therefore, in order to retain the accuracy level of the specified problem, this computation of the \vec{p}_ν must be done with special provisions. Although floating-point techniques are available for maintaining the accuracy in polynomial shifts under cancellation, it appears simplest to perform the computation of the \vec{p}_ν in *rational arithmetic*, with a subsequent return to floating-point data. Algebraic manipulations with polynomials will be needed throughout the following in any case, and every computer algebra system has rational arithmetic readily available.

We can now formulate our task in its final form. We will *drop the arrows* in the following for a simpler notation but retain ξ and ζ for the shifted variables and zeros.

Cluster analysis problem: Given a system P of s polynomials $p_\nu \in \mathbb{P}^s$, $\nu = 1(1)s$, with a cluster of m zeros about the origin, and a system $P_0 = \{p_{0\nu}, \nu = 1(1)s\}$, with an m -fold zero at the origin, with $P - P_0$ small (in a sense to be defined rigorously in the following). Find sufficiently good approximations for the individual zeros ζ_μ of P , $\mu = 1(1)m$, in the cluster.

2 The Univariate Case

For the case $s = 1$, the cluster analysis problem has been solved in [3]; see also the Ph.D. Thesis [2] of V. Hribernik. However, the presentation in [3] is not well suited as a basis for a generalization to $s > 1$. Therefore, in this section, we sketch the procedure for $s = 1$ in a form which may serve as a guideline for the multivariate procedure.

At first, we introduce a shorthand notation for polynomials of degree $m-1$; the usefulness of this notation will become obvious in the multivariate case:

$$\sum_{\mu=1}^m a_\mu \xi^{\mu-1} = (a_1 \dots a_m)^T \begin{pmatrix} 1 \\ \xi \\ \vdots \\ \xi^{m-1} \end{pmatrix} =: a^T \mathbf{t}, \quad a \in \mathbb{C}^m, \mathbf{t} := (1 \ \xi \ \dots \ \xi^{m-1})^T.$$

A polynomial p with an m -cluster at the origin may be written in the form

$$p(\xi) = (e^T \mathbf{t}) + (d_1^T \mathbf{t}) \xi^m + (d_2^T \mathbf{t}) \xi^{2m} + \dots, \quad (1)$$

with $\|e^T\|$ small and $|d_{11}|$ *not* small; a neighboring polynomial p_0 with a genuine m -fold zero at the origin is

$$p_0(\xi) = (d_1^T \mathbf{t}) \xi^m + (d_2^T \mathbf{t}) \xi^{2m} + \dots.$$

The reason for writing p and p_0 as “polynomials” in ξ^m , with coefficients from the span of the normal set of $\langle \xi^m \rangle$, will become clear below. The *distance* between p and p_0 may be defined by $\|e^T\|$.

Now consider an analogous representation of p as a polynomial in $(\xi^m + c^T \mathbf{t})$, with $c \in \mathbb{C}^m$:

$$p(\xi) = (e^T(c) \mathbf{t}) + (d_1^T(c) \mathbf{t}) (\xi^m + c^T \mathbf{t}) + (d_2^T(c) \mathbf{t}) (\xi^m + c^T \mathbf{t})^2 + \dots ; \quad (2)$$

The coefficients $e(c), d_1(c)$ etc. of this representation are uniquely defined as the remainders of a recursive division of p and its quotients by $(\xi^m + c^T \mathbf{t})$. (2) defines a map $F : \mathbb{C}^m \rightarrow \mathbb{C}^m$, which maps a coefficient vector $c^T \in \mathbb{C}^m$ into the residual vector $e^T(c)$ of the representation (2); this map satisfies

$$F(0) = e^T(0) = e^T .$$

Since $\|e^T\|$ is small, we may expect to find a small coefficient vector $c^{*T} \in \mathbb{C}^m$ such that

$$F(c^{*T}) = 0 . \quad (3)$$

Then

$$p(\xi) = (d_1^T(c^*) \mathbf{t}) (\xi^m + c^{*T} \mathbf{t}) + (d_2^T(c^*) \mathbf{t}) (\xi^m + c^{*T} \mathbf{t})^2 + \dots ,$$

and the m zeros of $g^*(\xi) := (\xi^m + c^{*T} \mathbf{t})$ are the individual zeros ζ_μ of the cluster. Note that g^* is simply the factor $\prod_{\mu=1}^m (\xi - \zeta_\mu)$ of p ; thus the existence of a solution c^{*T} of (3), with $\|c^{*T}\|$ small, is immediate.

Naturally, we do not strive to solve (3) exactly; instead we perform one Newton step, from $c^T = 0$: Since $F(c^T) = e^T(c)$, we obtain an approximation \tilde{c}^T for c^{*T} from

$$\tilde{c}^T F'(0) = -e^T . \quad (4)$$

To find the $m \times m$ -matrix $F'(0)$, we differentiate (2) w.r.t. c^T :

$$\mathcal{O} = F'(c^T) \mathbf{t} + (d_1^T(c^T) \mathbf{t}) I \mathbf{t} + \dots ,$$

where all further terms contain at least one factor $(\xi^m + c^T \mathbf{t})$. Setting $c^T = 0$ and observing that

$$(d_1^T \mathbf{t}) I \mathbf{t} \equiv \begin{pmatrix} d_{11} & d_{12} & \dots & d_{1m} \\ & d_{11} & \dots & d_{1,m-1} \\ & & \ddots & \vdots \\ & 0 & & d_{11} \end{pmatrix} \mathbf{t} \pmod{\xi^m} ,$$

we obtain

$$F'(0) = - \begin{pmatrix} d_{11} & d_{12} & \dots & d_{1m} \\ & d_{11} & \dots & d_{1,m-1} \\ & & \ddots & \vdots \\ & 0 & & d_{11} \end{pmatrix} .$$

Thus, our approximate coefficient vector \tilde{c}^T is obtained from the triangular linear system (4), with $d_{1\mu}$ and e_μ are from the original expansion (1) of p in powers of ξ^m .

Since the m zeros $\tilde{\zeta}_\mu$ of $\tilde{g}(\xi) = \xi^m + \tilde{c}^T \mathbf{t}$ are generally well isolated in a relative sense, any state-of-the-art (floating-point) polynomial solver will produce good approximations for them. For the same reason, individual improvement by Newton steps, starting from these

$\tilde{\zeta}_\mu$, will generally work. An exception is the occurrence of (secondary) clustering in \tilde{g} which could be remedied by another application of our procedure, presumably with a much smaller value of m . Thus the construction of \tilde{g} by the solution of a linear system of size m is the appropriate instrument for the “splitting” of the m -fold zero $\zeta = 0$ into approximations $\tilde{\zeta}_\mu$ for the m individual cluster zeros ζ_μ .

For use in the multivariate case, we may interpret the Newton step for (3) as a *first order perturbation* procedure: We perturb the factor polynomial $g_0 = \xi^m$ of p_0 by a small element from the span of the associated normal set such that the perturbed polynomial $\tilde{g} = g_0 + c^T \mathbf{t}$ becomes a first order correct approximation of g^* :

$$\tilde{g} \equiv O(\|c\|^2) .$$

Altogether, this suggests the following fully algebraic interpretation of our approach as a basis for a multivariate generalization:

$\langle \xi^m \rangle$	primary ideal \mathcal{I}_0 of the m -fold zero at the origin
$g_0(\xi) = \xi^m$	Groebner basis \mathcal{G}_0 of \mathcal{I}_0
$I, D, \dots, D^{m-1} _{\xi=0}$	linear functionals which form a dual basis \mathcal{D}_0 of \mathcal{I}_0 ($D := \partial / \partial \xi$)
$\{1, \xi, \dots, \xi^{m-1}\}$	normal set \mathcal{N} of \mathcal{G}_0 , basis of the residue class ring $\mathcal{R}_0 = \mathbb{P} / \mathcal{I}_0$
$\mathbf{t} = (1, \dots, \xi^{m-1})^T$	vector of (ordered) terms in \mathcal{N}
$p = (e^T \mathbf{t}) + \sum (d_j^T \mathbf{t})(g_0)^j$	representation of p in terms of \mathcal{G}_0
$\xi^m + c^t \mathbf{t}$	perturbation of \mathcal{G}_0 by an element from $\overline{\mathcal{N}} = \text{span } \mathcal{N}$
$g^*(\xi) = \xi^m + c^{*T} \mathbf{t}$	basis \mathcal{G}^* of the primary ideal \mathcal{I}^* of the m -cluster
$\tilde{g}(\xi) = \xi^m + \tilde{c}^T \mathbf{t}$	1st order correct approximation of g^*
$\tilde{g}(\xi) = 0$	generates approximations $\tilde{\zeta}_\mu$ of the m cluster zeros, i.e. it splits the m -fold zero

3 Multivariate Cluster Analysis

3.1 Layout of the Procedure

In accordance with the problem formulation at the end of section 1, we begin with systems $P_0 = \{p_{0\nu}, \nu = 1(1)s\}$ and $P = \{p_\nu, \nu = 1(1)s\}$ of polynomials from \mathbb{P}^s . P_0 has an m -fold zero at the origin; therefore, the neighboring system P must have m zeros close to the origin since we have assumed a *complete intersection* situation where zeros depend continuously on the data and cannot disappear except to infinity. We will now generalize our univariate approach described in section 2, cf. the table at the end of that section.

In a *first stage* of our analysis, we must find a *dual basis* \mathcal{D}_0 and a Groebner basis \mathcal{G}_0 of the primary ideal \mathcal{I}_0 of the m -fold zero of P_0 . This is no longer trivial: The *differential structure* of an m -fold zero in \mathbb{C}^s may take many different forms. Readers unfamiliar with this situation

are referred to [5] and [7], e.g., for a detailed analysis. The backbone which connects \mathcal{D}_0 and \mathcal{G}_0 is the residue class ring $\mathcal{R}_0 = \mathbb{P}^s / \mathcal{I}_0$ and its basis \mathcal{N} , the normal set of \mathcal{I}_0 . Naturally, the appearance of these objects depends on the chosen *term order*; in the following we assume that a consistent term order has been specified and is kept fixed throughout. The following sequence of steps appears appropriate:

1) From P_0 , find a basis $\{L_\mu, \mu = 1(1)m\}$ of the closed vector space U of differentials which describes the structure of the m -fold zero of P_0 at the origin. These linear functionals form a dual basis \mathcal{D}_0 of \mathcal{I}_0 ; cf. [5].

2) From \mathcal{D}_0 , determine a closed set \mathcal{N} of terms $t_\mu, \mu = 1(1)m$, such that the *Gram matrix* $(L_\mu[\mathbf{t}])$, with $\mathbf{t} := (t_1, \dots, t_m)^T$ is nonsingular. \mathcal{N} is a basis of \mathcal{R}_0 .

3) From \mathcal{D}_0 and \mathcal{N} , determine the multiplication table matrices $A_\sigma, \sigma = 1(1)s$, which describe the multiplicative structure of \mathcal{R}_0 . From these A_σ , the elements g_κ of the Groebner basis \mathcal{G}_0 are obtained immediately.

So far we have not dealt with the given system P at all. In the *second stage* of our analysis, we represent P in terms of \mathcal{G}_0 ; the residual vectors e_ν^T in $\text{NF}_{\langle \mathcal{G}_0 \rangle} p_\nu = e_\nu^T \mathbf{t}$ display the distance of P from P_0 . If this distance is *significant* (cf. section 1) we must split the m -fold zero by perturbing \mathcal{G}_0 into an approximation $\tilde{\mathcal{G}}_\varepsilon$ for the basis \mathcal{G}^* of the primary ideal \mathcal{I}^* of the cluster. This stage consists of the following steps:

4) Expand the p_ν of $P, \nu = 1(1)s$, into “polynomials” in the elements g_κ of \mathcal{G}_0 :

$$p_\nu = (e_\nu^T \mathbf{t}) + \sum_{\kappa=1}^k (d_{\nu\kappa}^T \mathbf{t}) g_\kappa + \sum_{\kappa_1 \geq \kappa_2}^k (d_{\nu\kappa_1\kappa_2}^T \mathbf{t}) g_{\kappa_1} g_{\kappa_2} + \dots \quad (5)$$

If $k = s$, the expansion is unique (see section 3.3); for $k > s$, we must also find and expand a generating set of *syzygies* of \mathcal{G}_0 .

5) Modify each element $g_\kappa \in \mathcal{G}_0$ by a generic *perturbation* $c_\kappa^T \mathbf{t}$ and consider the analogous expansions of the p_ν which define residual vectors $e_\nu^T(c)$, with $e_\nu^T(0) = e_\nu^T$ of (5); find $c^{*T} = (c_1^{*T}, \dots, c_k^{*T})$ such that $e_\nu^T(c^*) = 0, \nu = 1(1)s$. For this purpose, we *linearize* the relation between the c_κ^T and the e_ν^T and compute the matrices which represent the linearized relations. If $k > s$, we also have to perturb the syzygy equations of \mathcal{G}_0 and linearize them w.r.t. to the perturbation coefficients c_κ^T .

6) Solve the *linear system of equations* obtained in step 5, with the original residuals e_ν^T from (5) as right-hand sides; this yields perturbation coefficients \tilde{c}_κ^T which are first order correct approximations of the c_κ^{*T} .

In the *third* and last stage of our procedure, we compute approximations $\tilde{\zeta}_\mu$ for the m zeros $\zeta_\mu, \mu = 1(1)m$, of \mathcal{I}^* which are the zeros of P in its m -cluster about the origin:

7) Compute approximations \tilde{A}_σ for the multiplication tables A_σ^* of $\mathcal{R}^* = \mathbb{P}^s / \mathcal{I}^*$, from the $\tilde{g}_\kappa = g_\kappa + \tilde{c}_\kappa^T \mathbf{t}$ obtained in step 6, cf. section 3.4. Here we retain the normal set \mathcal{N} as basis for the representation of \mathcal{R}^* .

8) Compute the eigenvectors of the \tilde{A}_σ by any state-of-the-art matrix eigenproblem code (e.g. from [8]). In the rare case of secondary clustering, such a cluster may be resolved by the same procedure that we have just described. The (approximately) joint eigenvectors

(normalized with first component 1) of the \tilde{A}_σ are interpreted as evaluations of the normal set vector \mathbf{t} at the $\tilde{\zeta}_\mu$.

9) Check the residuals $p_\nu(\tilde{\zeta}_\mu)$ – evaluated with appropriate care – for significance relative to our specified accuracy level. If some $\tilde{\zeta}_\mu$ has significant residuals we perform a classical Newton step from $\tilde{\zeta}_\mu$ for further refinement. Otherwise, we accept the $\tilde{\zeta}_\mu$, $\mu = 1(1)m$, as sufficiently good approximations for the individual zeros in the cluster.

The complete procedure may appear unduly complicated and lengthy. But all numerical computations may be performed in floating-point arithmetic and, generally, m will not be a very large number; thus the computational effort will be small. In any case, we do not know of another algorithmic procedure which splits an m -fold zero of a multivariate polynomial system with high accuracy. In the following subsections, we will explain the individual steps with more detail.

3.2 Determination of the Primary Ideal of the m -fold Zero

In \mathbb{C}^s , an m -fold zero z_0 of a set P_0 of polynomials from \mathbb{P}^s is characterized by an m -dimensional vector space U of linear differential operators which satisfy

$$L[p](z_0) = 0, \quad \forall L \in U, p \in P_0. \quad (6)$$

The set of *all* polynomials for which (6) holds is an *ideal* if and only if

$$L[\cdot] \in U \quad \Rightarrow \quad L[x_\sigma \cdot] \in U, \quad \sigma = 1(1)s, \quad (7)$$

see, e.g., [6]; then U is called *closed*. Thus step 1 of our procedure requires the construction of a basis \mathcal{D}_0 for the closed vector space U of linear differential operators which characterize the m -fold zero $\zeta_0 = 0$ of our set P_0 . \mathcal{D}_0 is then a *dual basis* for the primary ideal \mathcal{I}_0 of this m -fold zero of P_0 .

This task has been discussed in [7] where an algorithm for its solution has also been proposed. It is based on the fact that a criterion for the closedness of U is (cf. [6])

$$L \in U \quad \Rightarrow \quad \sigma_i L \in U, \quad i = 1(1)s, \quad (8)$$

where σ_i denotes “antidifferentiation” w.r.t. x_i :

$$\text{For } L = D(x^j) := \frac{1}{j_1! \dots j_s!} \partial_{x_1^{j_1}} \dots \partial_{x_s^{j_s}}, \quad \sigma_i L := D(x^j/x_i), \quad (9)$$

where $D(x^j/x_i) = \mathcal{O}$ (zero functional) if $x_i \nmid x^j$; furthermore, $\sigma_i(\sum c_\mu L_\mu) = \sum c_\mu \sigma_i L_\mu$.

The construction of \mathcal{D}_0 therefore begins with $L_1 = D(1) = I$ and proceeds by checking suitable candidates with leading differential $D(x^j)$ in increasing term order: These are tested for consistency with the closedness criterion (8) and against the duality condition (6) for P_0 . The candidates are linear combinations of derivatives of functionals already in \mathcal{D}_0 and of functionals which have previously been found consistent with (8) but not (6); if coefficients for the linear combination exist such that (8) and (6) are satisfied, a new element of \mathcal{D}_0 has

been found. The algorithm stops when no further leading differentials consistent with (7) are available. This also determines m if it was not known a priori.

According to the authors of [7], their algorithm has never been implemented. In his complete implementation of the cluster analysis procedure, G. Thallinger has also implemented the above construction procedure for \mathcal{D}_0 in a slightly modified form. The algorithm and its implementation are nontrivial; details will be published in [9].

After the dual basis $\mathcal{D}_0 = \{L_1, \dots, L_m\}$ of \mathcal{I}_0 has thus been determined, we form the normal set \mathcal{N} for a representation of \mathcal{I}_0 w.r.t. the chosen term order in step 2 of our procedure: We consider the terms t of \mathbb{P}^s in increasing order, beginning with $t_1 = 1$, and form the row vectors $L[t] := (L_1[t], \dots, L_m[t])$, with all functionals evaluated at $\zeta_0 = 0$. If $L[t]$ is *linearly independent* of the vectors $L[t_\mu]$ for the t_μ already in \mathcal{N} , we admit t into \mathcal{N} , too. If it is not, we *delete its multiples* from further consideration. After m elements have been selected, the set $\mathcal{N} = \{t_1, \dots, t_m\}$ which we have obtained satisfies

$$t_\mu = x^j \in \mathcal{N} \quad \Rightarrow \quad x^j / x_i \in \mathcal{N} \quad \text{for } x_i \mid x^j, \quad i = 1(1)s, \quad (10)$$

i.e. \mathcal{N} is closed. The closedness of \mathcal{N} is a necessary and sufficient condition that it can function as a normal set of an ideal. From now on, $t_\mu, \mu = 1(1)m$, will denote the terms in \mathcal{N} and \mathbf{t} the (column) vector $(t_1 \dots t_m)^T$.

In step 3 of our procedure, we interpret \mathcal{N} as basis of the residue class ring $\mathcal{R}_0 = \mathbb{P}^s / \mathcal{I}_0$ and determine the multiplication table matrices A_σ of its multiplicative structure. We observe that

$$x_\sigma \mathbf{t} \equiv A_\sigma \mathbf{t} \pmod{\mathcal{I}_0} \quad \Leftrightarrow \quad L_\mu[x_\sigma \mathbf{t}] = A_\sigma L_\mu[\mathbf{t}], \quad L_\mu \in \mathcal{D}_0,$$

since \mathcal{D}_0 is a dual basis of \mathcal{I}_0 . Hence

$$A_\sigma = L[x_\sigma \mathbf{t}] \cdot L[\mathbf{t}]^{-1}, \quad \sigma = 1(1)s, \quad (11)$$

where $L[\mathbf{t}]$ and $L[x_\sigma \mathbf{t}]$ are the matrices formed by the row vectors $L_\mu[\mathbf{t}]$ and $L_\mu[x_\sigma \mathbf{t}]$, resp.

Remark: Algorithmically, the checking of the linear independence of the $L[t_\mu]$ generates $L[\mathbf{t}]$ in the *triangular* form $T = M_m \dots M_2 L[\mathbf{t}] P_1 \dots P_{m-1}$, with elimination matrices M_μ and permutations P_μ . Hence $L[\mathbf{t}]^{-1} = P_1 \dots P_{m-1} T^{-1} M_m \dots M_2$ in the evaluation of (11). Note that $L[\mathbf{t}]$ is the Gram matrix of the *interpolation problem*

$$\text{Find } p \in \text{span } \mathcal{N} : \quad L_\mu[p]_{\zeta_0=0} = l_\mu \in \mathbb{C}^s, \quad \mu = 1(1)m. \quad \square \quad (12)$$

From \mathcal{N} , we obtain a set $\text{LT}[\mathcal{G}]$ of leading terms for a generating set \mathcal{G}_0 of \mathcal{I}_0 :

$$\text{LT}[\mathcal{G}] := \{t : t \notin \mathcal{N}, t/x_i \in \mathcal{N} \text{ for } x_i \mid t, i = 1(1)s\}. \quad (13)$$

Obviously, each element of $\text{LT}[\mathcal{G}]$ is a component of one or several vector(s) $x_\sigma \mathbf{t}$; hence the corresponding row of $A_\sigma \mathbf{t}$ defines the remaining terms of the respective element g_κ of \mathcal{G} .

Proposition 3.1: The generating set \mathcal{G} thus obtained is the Groebner basis \mathcal{G}_0 of \mathcal{I}_0 for the specified term order.

Proof: Let \mathcal{N}_0 be the normal set of \mathcal{G}_0 . If $\mathcal{N}_0 = \mathcal{N}$, the polynomials g_κ from the A_σ must be the elements of \mathcal{G}_0 . Assume $\mathcal{N}_0 \neq \mathcal{N}$; since $|\mathcal{N}_0| = |\mathcal{N}|$, there exists $\bar{t} := \min$

$\{t \in \mathcal{N}, t \notin \mathcal{N}_0\}$. By (10), for all σ with $x_\sigma | \bar{t}$, $\bar{t}/x_\sigma \in \mathcal{N}$ and (trivially) $\bar{t}/x_\sigma < \bar{t}$, hence $\bar{t}/x_\sigma \in \mathcal{N}_0$. Since $\bar{t} \notin \mathcal{N}_0$ this implies $\bar{t} \in \text{LT}[\mathcal{G}_0]$; the remaining terms in the respective Groebner basis element g must be smaller than \bar{t} and from \mathcal{N}_0 , hence also from \mathcal{N} . This implies $g \in \text{span } \mathcal{N}$ which is a contradiction. \square

Remark: While our determination of \mathcal{N} naturally follows the respective algorithm in [6], our computation of the Groebner basis via (11) differs slightly from the approach in [6]. \square

3.3 Perturbation of \mathcal{I}_0 into the Primary Ideal of the Cluster

We have now obtained a complete quantitative representation of the primary ideal \mathcal{I}_0 of the m -fold zero $\zeta_0 = 0$ of the polynomial set P_0 through its Groebner basis \mathcal{G}_0 , its dual basis \mathcal{D}_0 , and its residue class ring \mathcal{R}_0 with basis \mathcal{N} and the multiplication tables A_σ .

These quantities now serve as our reference in the analysis of the given set P of polynomials which has an m -cluster of zeros $\zeta_\mu, \mu = 1(1)m$, about the origin. We can immediately compute the *residual vectors* $e_\nu^T \in \mathbb{C}^m$ of the p_ν , from

$$\text{NF}_{\langle \mathcal{G}_0 \rangle} p_\nu = \sum_{\mu=1}^m e_{\nu\mu} t_\mu := e_\nu^T \mathbf{t} \quad (14)$$

by a normal form algorithm; they are unique since \mathcal{G}_0 is a Groebner basis.

Naturally, we expect the residual vectors e_ν^T to be small since we have assumed P to be close to P_0 . On the other hand, for the following to be meaningful for *not fully accurate data* (cf. case 2 in section 1), the size of the e_ν^T must be *significant* relative to the specified data accuracy: If the polynomials $\tilde{p}_\nu := p_\nu - e_\nu^T \mathbf{t} \in \mathcal{I}_0$ constitute a set \tilde{P} in the equivalence class \overline{P} , then this equivalence class contains polynomial systems which have a genuine m -fold zero at the origin and thus the origin is an m -fold *pseudozero* of P . We will assume that this is *not* the case throughout the following.

For the further analysis, we also need the coefficient vectors $d_{\nu\kappa}^T$ of (5). They can be computed from representations

$$p_\nu - e_\nu^T \mathbf{t} = \sum_{\kappa=1}^k q_{\nu\kappa} g_\kappa, \quad \nu = 1(1)s, \quad (15)$$

by further application of the NF-algorithm:

$$d_{\nu\kappa}^T \mathbf{t} = \text{NF}_{\langle \mathcal{G}_0 \rangle} q_{\nu\kappa}, \quad \kappa = 1(1)k, \quad \nu = 1(1)s. \quad (16)$$

The representations (15) are unique except for *syzygies*:

$$\sum_{\kappa=1}^k s_{\lambda\kappa} g_\kappa = \mathcal{O} \quad (\text{the zero polynomial}), \quad \lambda = 1(1)l, \quad (17)$$

where the λ refer to some generating set of the syzygy module \mathcal{S}_0 of \mathcal{G}_0 .

In the case $k = s$, all $\text{LT}[g_\kappa]$ are monomials $\xi_\sigma^{j_\sigma}, \sigma = 1(1)s$; thus a generating set of \mathcal{S}_0 consists only of the *trivial* syzygies $g_{\kappa_1} g_{\kappa_2} - g_{\kappa_2} g_{\kappa_1} = \mathcal{O}$, and such syzygies do not affect the $d_{\nu\kappa}^T$ in (16). Hence the $d_{\nu\kappa}^T$ – and all further coefficient vectors in (5) – are *unique* for $k = s$.

Unfortunately, this case of a *rectangular* normal set rarely occurs in our primary ideals of m -fold zeros. For $k > s$, we denote the $k - s$ elements in \mathcal{G}_0 whose leading term is *not* a monomial by g_{s+1}, \dots, g_k . The elements of a generating set of \mathcal{S}_0 are then formed as follows: Let $\text{LT}[g_\kappa], \kappa \in \{s+1, \dots, k\}$, contain a power of x_{σ_κ} ; then form the S-polynomial $S[g_\kappa, g_{\sigma_\kappa}]$ where the leading term of g_{σ_κ} is a power of x_{σ_κ} and represent it in terms of \mathcal{G}_0 . The collection of all identities obtainable in this fashion is a generating set of \mathcal{S}_0 (cf., e.g., [10]).

In analogy to our transition from (15) to (16), we now form coefficient vectors $s_{\lambda\kappa}^T$ through

$$s_{\lambda\kappa}^T \mathbf{t} := \text{NF}_{\langle \mathcal{G}_0 \rangle} s_{\lambda\kappa}, \quad \kappa = 1(1)k, \quad \lambda = 1(1)l, \quad (18)$$

where the polynomials $s_{\lambda\kappa}$ on the right-hand side are from (17). These $s_{\lambda\kappa}^T \in \mathbb{C}^m$ characterize the ambiguity in the $d_{\nu\kappa}^T$ of (16).

We have now finished step 4 in the second stage of our procedure. In step 5, we have to consider the mapping $F : \mathbb{C}^{k \cdot m} \rightarrow \mathbb{C}^{s \cdot m}$ which is defined by expansions like (5) of the p_ν in terms of *perturbed Groebner basis elements* $g_\kappa + c_\kappa^T \mathbf{t}, \kappa = 1(1)k$:

$$p_\nu = (e_\nu^T(c) \mathbf{t}) + \sum_{\kappa=1}^k (d_{\nu\kappa}^T(c) \mathbf{t}) (g_\kappa + c_\kappa^T \mathbf{t}) + \sum_{\kappa_1 \geq \kappa_2}^k (d_{\nu\kappa_1\kappa_2}^T(c) \mathbf{t}) (g_{\kappa_1} + c_{\kappa_1}^T \mathbf{t}) (g_{\kappa_2} + c_{\kappa_2}^T \mathbf{t}) + \dots \quad (19)$$

through $F(c_\kappa^T, \kappa = 1(1)k) := (e_\nu^T(c), \nu = 1(1)s)$.

Such *Perturbed Groebner Basis Representations* have been considered in our paper [4]; there it has been shown that they are the appropriate means to represent all polynomials in a neighborhood of a degenerate set of polynomials in a continuous and uniform way. We have assembled some basic information about perturbed Groebner bases in the Appendix.

In the case $k > s$, where there are more perturbation vectors c_κ^T than residual vectors e_ν^T , the representation (19) has to be supplemented by the analogously perturbed version of (17):

$$\mathcal{O} = \sum_{\kappa=1}^k (s_{\lambda\kappa}^T(c) \mathbf{t}) (g_\kappa + c_\kappa^T \mathbf{t}) + \sum_{\kappa_1 \geq \kappa_2}^k (s_{\lambda\kappa_1\kappa_2}^T(c) \mathbf{t}) (g_{\kappa_1} + c_{\kappa_1}^T \mathbf{t}) (g_{\kappa_2} + c_{\kappa_2}^T \mathbf{t}) + \dots, \quad (20)$$

$\lambda = 1(1)l$, which provides the restrictions on the c_κ^T needed to make the *restricted mapping* $\hat{F} : \mathbb{C}^{k \cdot m} \cap (20) \rightarrow \mathbb{C}^{s \cdot m}$ *bijective* in a neighborhood of $(c_\kappa^T) = 0$; cf. the Appendix.

Obviously, the solution $c^* = (c_1^{*T}, \dots, c_k^{*T})$ of the polynomial system in the c_κ

$$\hat{F}(c_1^T, \dots, c_k^T) = (0, \dots, 0) \quad (21)$$

generates the primary ideal $\mathcal{I}^* = \langle \mathcal{G}^* \rangle := \langle g_1 + c_1^{*T} \mathbf{t}, \dots, g_k + c_k^{*T} \mathbf{t} \rangle$ of the m -cluster of P . However, the exact determination of the c_κ^{*T} is far too involved; instead, we determine perturbation vectors \tilde{c}_κ^T which solve the *linearization of (21) about $c_\kappa^T = 0$* . This simplification of the task posed by (21) may also be considered as *one Newton step* for (21) from $c = 0$; the approximate elements $\tilde{g}_\kappa := g_\kappa + \tilde{c}_\kappa^T \mathbf{t}$ satisfy

$$\tilde{g}_\kappa - g_\kappa^* = O(\|\tilde{c}\|^2), \quad \kappa = 1(1)k. \quad (22)$$

To determine the *linear system of equations* which defines these \tilde{c}_κ^T , we set

$$e_\nu^T(c) = e_\nu^T + \sum_{\kappa=1}^k c_\kappa^T E_{\nu\kappa} + O(\|c_\kappa^T\|^2), \quad (23)$$

$$d_{\nu\kappa}^T(c) = d_{\nu\kappa}^T + O(\|c_\kappa^T\|), \quad (24)$$

and, in the case $k > s$, analogously

$$s_{\lambda\kappa}^T(c) = s_{\lambda\kappa}^T + O(\|c_\kappa^T\|), \quad (25)$$

for the respective values of the subscripts on the left-hand side. Then we introduce (23), (24), and (25) into (19) and (20) resp. and subtract the unperturbed representations (5) and (17) resp. Since the left-hand sides vanish, the normal form of the right-hand sides w.r.t. \mathcal{I}_0 must vanish. Thus we obtain, for $\nu = 1(1)s$,

$$\mathcal{O} = \sum_{\kappa=1}^k (c_\kappa^T E_{\nu\kappa}) \mathbf{t} + \text{NF}_{\langle \mathcal{G}_0 \rangle} \sum_{\kappa=1}^k (d_{\nu\kappa}^T \mathbf{t})(c_\kappa^T \mathbf{t}) + O(\|c\|^2) \quad (26)$$

and, for $\lambda = 1(1)l$ in the case $k > s$,

$$\mathcal{O} = \text{NF}_{\langle \mathcal{G}_0 \rangle} \sum_{\kappa=1}^k (s_{\lambda\kappa}^T \mathbf{t})(c_\kappa^T \mathbf{t}) + O(\|c\|^2); \quad (27)$$

all terms not explicitly spelt out above either contain at least one factor $g_\kappa \in \mathcal{G}_0$ or they are quadratic in the c_κ^T (or both). The normal forms of products of elements from $\text{span } \mathcal{N}$ are computed by

Proposition 3.2: Denote the multiplication table matrices for $t_\nu = x^{j(\nu)} \in \mathcal{N}$ by $A_{t_\nu} := \prod A_1^{j_1^{(\nu)}} \dots A_s^{j_s^{(\nu)}}$. Then

$$\text{NF}_{\mathcal{I}_0}(a^t \mathbf{t})(b^T \mathbf{t}) = [b^T(a^T \mathbf{t}(A))] \mathbf{t},$$

where

$$\mathbf{t}(A) := (A_{t_1} \dots A_{t_m})^T, \quad (a^T \mathbf{t}(A)) := \sum_{\nu=1}^m a_\nu A_{t_\nu}.$$

Proof:

$$\begin{aligned} (a^t \mathbf{t})(b^T \mathbf{t}) &= \left(\sum_{\nu} a_\nu t_\nu \right) \left(\sum_{\mu} b_\mu t_\mu \right) = \sum_{\mu} b_\mu \left(\sum_{\nu} a_\nu t_\nu t_\mu \right) \\ &\equiv_{\mathcal{I}_0} \sum_{\mu} b_\mu \left(\sum_{\nu} a_\nu A_{t_\nu} \right)_\mu \mathbf{t} = \sum_{\mu} b_\mu (a^T \mathbf{t}(A))_\mu \mathbf{t} \end{aligned}$$

where $(\text{matrix})_\mu$ denotes the μ -th row of the matrix. \square

Thus, (26) and (27) imply

$$-E_{\nu\kappa} = (d_{\nu\kappa}^T \mathbf{t}(A)) \quad \text{and} \quad S_{\lambda\kappa} := (s_{\lambda\kappa}^T \mathbf{t}(A)). \quad (28)$$

With (28), our request for $e_\nu^T(c) = O(\|c_\kappa^T\|^2)$ in (23) turns into the following linear system for the c_κ^T :

$$(c_1^T \dots c_k^T) \begin{pmatrix} -E_{11} & \dots & -E_{s1} & S_{11} & \dots & S_{l1} \\ \vdots & & \vdots & \vdots & & \vdots \\ -E_{1k} & \dots & -E_{sk} & S_{1k} & \dots & S_{lk} \end{pmatrix} = (e_1^T \dots e_s^T 0 \dots 0). \quad (29)$$

The system (29) looks overdetermined for $l > k - s$. However,

$$\text{rk} \begin{pmatrix} S_{11} & \dots & S_{l1} \\ \vdots & & \vdots \\ S_{1k} & \dots & S_{lk} \end{pmatrix} = (k - s) m \quad (30)$$

and the rank of the complete matrix in (29) is $k m$ so that (29) has a unique solution \tilde{c}_κ^T , $\kappa = 1(1)k$. Details are explained in the Appendix.

By (22), the perturbation vectors \tilde{c}_κ^T from (29) provide good approximations for the elements g_κ^* of the perturbed Groebner basis \mathcal{G}^* of the primary ideal \mathcal{I}^* of the m zeros of P clustered about the origin.

3.4 Computation of the Cluster Zeros

The cluster zeros ζ_μ , $\mu = 1(1)m$, of P are the zeros of the polynomial system G^* :

$$g_\kappa^*(\xi) = g_\kappa(\xi) + c_\kappa^{*T} \mathbf{t}(\xi) = 0, \quad \kappa = 1(1)k,$$

whose zeros are normally well-separated in a relative sense. If there is secondary clustering within the cluster, the primary ideal for these zeros can be found (after a shift to the origin) by another application of our approach; therefore we do not pursue this possibility further.

However, to obtain approximations $\tilde{\zeta}_\mu$ for the ζ_μ we cannot – except in the case $k = s$ – replace the system G^* by the system \tilde{G} :

$$\tilde{g}_\kappa(\xi) = g_\kappa(\xi) + \tilde{c}_\kappa^T \mathbf{t}(\xi) = 0, \quad \kappa = 1(1)k,$$

because this system will generally be inconsistent and may have no zeros at all. Instead we must use the information from the \tilde{c}_κ^T to solve G^* approximately. For this purpose, we compute approximations \tilde{A}_σ for the multiplication table matrices A_σ^* of the residue class ring $\mathcal{R}^* = \mathbb{P}^s / \mathcal{I}^*$, with basis \mathcal{N} . Since the joint eigenvectors of the *commuting family* of $m \times m$ -matrices A_σ are the *evaluations* $\mathbf{t}(\zeta_\mu)$ of the normal set vector \mathbf{t} at the ζ_μ (after normalization of the first components to 1; cf., e.g., [11]), the eigenvectors of the \tilde{A}_σ will provide approximations for the $\mathbf{t}(\zeta_\mu)$ and thus for the ζ_μ .

At first, we remember that the rows $a_{\sigma\mu}^{*T}$, $\mu = 1(1)m$, of the A_σ^* are determined by $\xi_\sigma t_\mu \equiv a_{\sigma\mu}^{*T} \mathbf{t} \pmod{\mathcal{G}^*}$. Naturally, if $\xi_\sigma t_\mu = t_{\mu_\sigma} \in \mathcal{N}$, we set $\tilde{a}_{\sigma\mu}^T := a_{\sigma\mu}^{*T} = e_{\mu_\sigma}^T$. If $\xi_\sigma t_\mu = \text{LT}[g_\kappa]$, we set

$$\tilde{a}_{\sigma\mu}^T \mathbf{t} := -(\tilde{g}_\kappa - \xi_\sigma t_\mu) \approx -(g_\kappa^* - \xi_\sigma t_\mu).$$

Each remaining $\xi_\sigma t_\mu$ is in the *border set* of \mathcal{N} (cf. [6]) and thus a multiple of some $\text{LT}[g_\kappa]$. If, e.g., $\xi_\sigma t_\mu = \xi_{\sigma'} \text{LT}[g_\kappa] = \xi_{\sigma'}(\xi_{\sigma_\kappa} t_{\mu_\kappa})$, we set

$$\tilde{a}_{\sigma\mu}^T := \tilde{a}_{\sigma_\kappa\mu_\kappa}^T \tilde{A}_{\sigma'}$$

since this relation would hold for the A_σ^* . In this way, all $\tilde{a}_{\sigma\mu}^T$ can be found recursively or from linear relations between them.

In step 8, we compute the m eigenvectors $\tilde{v}_{\sigma\mu} \in \mathbb{C}^s$, $\mu = 1(1)m$, of the \tilde{A}_σ by some linear algebra package (e.g. [8]) and normalize their first components to 1 (these components cannot vanish; see [11]). Naturally, our approximate multiplication table matrices \tilde{A}_σ are not strictly commuting; some elements in $A_{\sigma_1}A_{\sigma_2} - A_{\sigma_2}A_{\sigma_1}$ may be of $O(\|\tilde{c}\|^2)$. Thus the eigenvector systems of the individual \tilde{A}_σ may differ slightly. However, due to $\tilde{A}_\sigma \approx A_\sigma^*$, we have

$$\tilde{v}_{\sigma\mu} \approx \mathbf{t}(\tilde{\zeta}_\mu), \quad \mu = 1(1)m, \quad (31)$$

Only if two or more of the $\tilde{\zeta}_\mu$ have extremely close values in their ξ_σ -components, the associated eigenvectors may not satisfy (31) individually but only span the respective eigenspace approximately. This situation may be diagnosed by consistency checks between the components of the $\tilde{v}_{\sigma\mu}$.

Since it is generally not clear which \tilde{A}_σ yields the best approximations for the ζ_μ , we have used the arithmetic means of the matching $\tilde{v}_{\sigma\mu}$ to obtain vectors \tilde{v}_μ which we interpret as $\mathbf{t}(\tilde{\zeta}_\mu)$. Each variable ξ_σ is either a term in \mathcal{N} or the leading term of a g_κ ; thus the components of the zeros ζ_μ may either be read from the \tilde{v}_μ directly or computed from the other components via a \tilde{g}_κ .

After we have obtained these approximations $\tilde{\zeta}_\mu$, $\mu = 1(1)m$, for all zeros in the m -cluster of P about the origin and thus successfully split the m -fold zero of P_0 , we must check whether a further numerical improvement of the $\tilde{\zeta}_\mu$ is necessary (Step 9): We compute the residuals $p_\nu(\tilde{\zeta}_\mu)$, preferably in rational arithmetic to cope with the heavy cancellation of leading digits.

To judge the significance of the residuals, it may suffice to consider their order of magnitude relative to the specified accuracy level. For a more refined judgement, we consider the *interpolation* of each residual vector $(p_\nu(\tilde{\zeta}_1) \dots p_\nu(\tilde{\zeta}_m))$ by a polynomial from $\overline{\mathcal{N}}$: Let $\rho_\nu^T \mathbf{t}(\tilde{\zeta}_\mu) = p_\nu(\tilde{\zeta}_\mu)$, $\mu = 1(1)m$, for $\nu = 1(1)s$; then the neighboring system $P_\rho = \{p_\nu - \rho_\nu^T \mathbf{t}, \nu = 1(1)s\}$ has exact zeros at the $\tilde{\zeta}_\mu$.

If all $\|\rho_\nu^T\|$ are below our specified accuracy limit, we can safely accept the $\tilde{\zeta}_\mu$ as sufficiently good approximations of the cluster zeros. Otherwise, we may improve a $\tilde{\zeta}_\mu$ which generates unduly large residuals by a standard Newton step from that $\tilde{\zeta}_\mu$. For arbitrarily accurate data, each zero can – in principle – be approximated individually to arbitrary accuracy since we are now close enough for the local convergence of the Newton method.

4 Example

Exclusively for reasons of display, we have used an indeterminate perturbation parameter ε in steps 1 through 5. Beginning with step 6, we have taken $\varepsilon = 10^{-4}$ and displayed only a

few digits of the resulting decimal fractions, because the rational expressions in ε would have been prohibitive for printing. Steps 8 and 9 require floating-point computation in any case.

In agreement with the task formulated at the end of section 1, we start with a system P_0 in 3 variables ξ_1, ξ_2, ξ_3 which has a multiple zero at the origin, and with a near-by system P whose zero cluster about the origin we wish to find:

$$P_0 = \begin{cases} p_{01}(\xi_1, \xi_2, \xi_3) &= \xi_1 + \xi_2^2 - 4\xi_2\xi_3 + 4\xi_3^2 \\ p_{02}(\xi_1, \xi_2, \xi_3) &= \xi_1^2 - 2\xi_1 + \xi_2^2 + \xi_3^2 \\ p_{03}(\xi_1, \xi_2, \xi_3) &= \xi_1\xi_2\xi_3 + \xi_2^2\xi_3 + \xi_2\xi_3^2 \end{cases}$$

$$P = \begin{cases} p_1(\xi_1, \xi_2, \xi_3) &= p_{01}(\xi_1, \xi_2, \xi_3) + \varepsilon (2\xi_1\xi_2 - \xi_1\xi_3 + 3\xi_2\xi_3 - 8) \\ p_2(\xi_1, \xi_2, \xi_3) &= p_{02}(\xi_1, \xi_2, \xi_3) + \varepsilon (-\xi_1 - 3\xi_2 + 5\xi_3 + 1) \\ p_3(\xi_1, \xi_2, \xi_3) &= p_{03}(\xi_1, \xi_2, \xi_3) + \varepsilon (\xi_1^3 + 3\xi_1^2\xi_3 + 3\xi_1\xi_3^2 + \xi_3^3) \end{cases}$$

Throughout the following, the term order is lexicographic, with $\xi_1 > \xi_2 > \xi_3$.

Step 1: Determination of \mathcal{D}_0

It is obvious that $D(1)$, $D(\xi_3)$, $D(\xi_2)$ are in \mathcal{D}_0 and also $D(\xi_2^2 - \frac{3}{4}\xi_2\xi_3 - \xi_3^2)$ and $D(\xi_1 + \frac{9}{4}\xi_2\xi_3 + 2\xi_3^2)$. Our algorithm finds no more 2nd order differential operators linearly independent of these and consistent with the closedness of U which vanish for P_0 . For a 3rd order differential operator which vanishes for p_{03} , the closedness criterion (8) is the main restriction. Our algorithm yields

$$L_6 = D(\xi_1\xi_2 - \frac{1}{21}\xi_1\xi_3 + \frac{17}{7}\xi_2^3 + \frac{3}{7}\xi_2^2\xi_3 - \frac{3}{7}\xi_2\xi_3^2 - \frac{11}{21}\xi_3^3), \quad \text{with}$$

$$\begin{aligned} \sigma_1 L_6 &= D(\xi_2 - \frac{1}{21}\xi_3) = -\frac{1}{21}L_2 + L_3, \\ \sigma_2 L_6 &= D(\xi_1 + \frac{17}{7}\xi_2^2 + \frac{3}{7}\xi_2\xi_3 - \frac{3}{7}\xi_3^2) = \frac{17}{7}L_4 + L_5, \\ \sigma_3 L_6 &= D(-\frac{1}{21}\xi_1 + \frac{3}{7}\xi_2^2 - \frac{3}{7}\xi_2\xi_3 - \frac{11}{21}\xi_3^2) = \frac{3}{7}L_4 - \frac{1}{21}L_5. \end{aligned}$$

No further closed extension of U is feasible; thus $m = 6$ and \mathcal{D}_0 consists of L_1, \dots, L_6 found above. Note the non-trivial form of L_6 in the specification of the differential structure of the 6-fold zero of P_0 at the origin.

Step 2: Determination of the normal set \mathcal{N} of \mathcal{I}_0

1 and the first three powers of ξ_3 yield independent row vectors $L[t_\mu]$, and so do ξ_2 and $\xi_2\xi_3$. Thus $\mathcal{N} = \{1, \xi_3, \xi_3^2, \xi_3^3, \xi_2, \xi_2\xi_3\}$, and

$$L[t] = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 2 & 0 \\ 0 & 0 & 0 & 0 & 0 & -\frac{11}{21} \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -\frac{3}{4} & \frac{9}{4} & 0 \end{pmatrix}.$$

Naturally, the linear algebra in the actual algorithm triangularizes $L[t]$ during its formation by permutations and eliminations; this information is used in

Step 3: Determination of the multiplication tables A_σ

$L[\xi_1 \mathbf{t}]$, $L[\xi_2 \mathbf{t}]$, and $L[\xi_3 \mathbf{t}]$ are easily evaluated and yield, via (11),

$$A_1 = \begin{pmatrix} 0 & 0 & -1 & 0 & 0 & \frac{4}{3} \\ 0 & 0 & 0 & \frac{1}{11} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -\frac{21}{11} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}, \quad A_2 = \begin{pmatrix} 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & \frac{9}{11} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -3 & 0 & 0 & \frac{8}{3} \\ 0 & 0 & 0 & -\frac{9}{11} & 0 & 0 \end{pmatrix},$$

$$A_3 = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & \frac{9}{11} & 0 & 0 \end{pmatrix}.$$

From (13), the $\text{LT}[g_\kappa]$ are $\{\xi_1, \xi_2^2, \xi_2\xi_3^2, \xi_3^4\}$ and the multiplication tables yield

$$\mathcal{G}_0 = \begin{cases} g_1 = \xi_1 - \frac{4}{3}\xi_2\xi_3 + \xi_3^2, \\ g_2 = \xi_2^2 - \frac{8}{3}\xi_2\xi_3 + 3\xi_3^2, \\ g_3 = \xi_3^4, \\ g_4 = \xi_2\xi_3^2 - \frac{9}{11}\xi_3^3. \end{cases}$$

as Groebner basis of the primary ideal \mathcal{I}_0 . In terms of \mathcal{G}_0 , P_0 appears as

$$p_{01} = g_1 + g_2,$$

$$p_{02} = (-2 - 2\xi_3^2 + \frac{8}{3}\xi_2\xi_3) g_1 + (1 + \frac{16}{9}\xi_3^2) g_2 + \frac{56}{27}\xi_3 g_4 - \frac{29}{11} g_3 + (g_1)^2,$$

$$p_{03} = \xi_2\xi_3 g_1 + (\xi_3 + \frac{4}{3}\xi_3^2 g_2 - \frac{21}{11} g_3 + (\frac{11}{3} + \frac{23}{9}\xi_3) g_4).$$

Step 4: Representation of the system P in terms of \mathcal{G}_0

By polynomial division, we obtain (cf. section 3.3):

$$p_1 = p_{01} + \varepsilon \left[(-8 - \frac{43}{11}\xi_3^3 + 3\xi_2\xi_3) + (-\xi_3 + 2\xi_2) g_1 + (\frac{8}{3}\xi_3) g_2 + (\frac{34}{9}) g_4 \right],$$

$$p_2 = p_{02} + \varepsilon \left[(1 + 5\xi_3 + \xi_3^2 - 3\xi_2 - \frac{4}{3}\xi_2\xi_3) + (-1) g_1 \right],$$

$$p_3 = p_{03} + \varepsilon \left[(\xi_3^3) + (3\xi_3^2 + \frac{6}{11}\xi_3^3) g_1 + (\frac{16}{3}\xi_3^3) g_2 + (\frac{3}{11} - \frac{87}{11}\xi_3 - \frac{431}{99}\xi_3^2) g_3 + (4\xi_3 + \frac{56}{9}\xi_3^2 - \frac{116}{243}\xi_3^3) g_4 + \dots \right],$$

where the dots indicate higher order terms in the g_κ . The representation of the syzygies of \mathcal{G}_0 takes the form (cf. (20))

$$\mathcal{O} = -\xi_3^2 g_2 + \frac{180}{121} g_3 + \left(-\frac{61}{33}\xi_3 + \xi_2\right) g_4, \quad \mathcal{O} = \left(-\frac{9}{11}\xi_3 + \xi_2\right) g_3 - \xi_3^2 g_4.$$

Step 5: Formation of the blocks of the matrix in (29)

The $E_{\nu\kappa}$ and $S\lambda\kappa$ are determined by (28), with $d_{\nu\kappa}^T$ and $s_{\lambda\kappa}^T$ from (5) and (17) resp. We indicate only the formation of E_{11} : From $d_{11}^T = (1, -1, 0, 0, 2, 0)$,

$$-E_{11} = I - \varepsilon A_3 + 2\varepsilon A_2 = \begin{pmatrix} 1 & -\varepsilon & 0 & 0 & 2\varepsilon & 0 \\ 0 & 1 & -\varepsilon & 0 & 0 & 2\varepsilon \\ 0 & 0 & 1 & \frac{7}{11}\varepsilon & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & -6\varepsilon & 0 & 1 & \frac{13}{3}\varepsilon \\ 0 & 0 & 0 & -\frac{27}{11}\varepsilon & 0 & 1 \end{pmatrix}.$$

All columns of the $S_{2\kappa}$ block column are linearly dependent on the 6 columns of the $S_{1\kappa}$ blocks.

Step 6: Solution of the linear system (29)

For $\varepsilon = 10^{-4}$, the final 24×24 system for the \tilde{c}_κ^T yields

$$\begin{aligned} \tilde{c}_1^T &= (-2.99990, -1.66644, +1.11105, -3.21078, +1.00017, -0.41657) \cdot 10^{-4}, \\ \tilde{c}_2^T &= (-5.00010, +1.66696, -1.11089, -0.69919, -0.99957, +3.41650) \cdot 10^{-4}, \\ \tilde{c}_3^T &= (+0.00000, +0.00000, -1.66670, +1.21714, +0.00000, -0.91669) \cdot 10^{-4}, \\ \tilde{c}_4^T &= (+0.00000, +1.36366, -0.45452, +0.70033, +0.00000, +0.61349) \cdot 10^{-4}. \end{aligned}$$

The perturbation coefficients are clearly of $O(\varepsilon)$. The $\tilde{g}_\kappa = g_\kappa + \tilde{c}_\kappa^T \mathbf{t}$ are now determined. Note that the coefficient of $\xi_2 \xi_3^2$ in g_3 is $\neq 0$; thus our basis \mathcal{G}_0 of the primary ideal of the cluster is *not* a genuine Groebner basis.

Step 7: Determination of the multiplication tables \tilde{A}_σ

From the \tilde{g}_κ , we obtain immediately (cf. section 3.4)

- the 1st row of \tilde{A}_1 ,
- the 3rd and 5th row of \tilde{A}_2 , the 1st and 2nd row are trivial,
- the 4th and 6th row of \tilde{A}_3 , all other rows are trivial.

The complete knowledge of \tilde{A}_3 permits the computation of $\tilde{a}_{24}^T := \tilde{a}_{23}^T \tilde{A}_3$ and $\tilde{a}_{26}^T := \tilde{a}_{25}^T \tilde{A}_3$ which completes \tilde{A}_2 . The remaining rows of \tilde{A}_1 are obtained analogously.

Step 8: Computation of the eigenvectors of the \tilde{A}_σ

The eigenvectors have been computed numerically by the respective Maple routine. Since the \tilde{A}_σ are not fully commutative (cf. section 3.4) their eigenvector systems differ slightly (after normalization of first components to 1). Also – due to the closeness of some eigenvalues – two eigenvectors of some \tilde{A}_σ are inconsistent as evaluations of \mathbf{t} at a point in \mathbb{C}^3 but span the correct eigenspace (approximately). We have not attempted to correct such eigenvectors but simply taken the arithmetic means of the matching eigenvectors.

The 2nd and 5th components of these “joint” eigenvectors are the ξ_3 - and ξ_2 -components of the approximate zeros $\tilde{\zeta}_\mu$, resp., $\mu = 1(1)6$. The ξ_1 -components are found from \tilde{g}_1 . This yields the following approximations for the 6 zeros in the cluster of P :

$$\begin{pmatrix} .0003022 \\ -.0223110 \\ 0 \end{pmatrix}, \begin{pmatrix} .0002977 \\ .0224109 \\ 0 \end{pmatrix}, \begin{pmatrix} .0001278 \\ -.0086180 \\ .0086604 \end{pmatrix}, \begin{pmatrix} .0001221 \\ .0087033 \\ -.0086612 \end{pmatrix}, \begin{pmatrix} .0001306 \\ -.0000965 \\ -.0129974 \end{pmatrix}, \begin{pmatrix} .0001360 \\ -.0001003 \\ .0128216 \end{pmatrix}.$$

Step 9: Residual formation, individual Newton steps

The evaluation of the p_ν at the above $\tilde{\zeta}_\mu$ yields residuals whose absolute values are all below $1.6 \cdot 10^{-6}$, many are considerably smaller. Also the components of the ρ_ν^T -vectors are below that value. This means that there exist polynomial systems with exact zeros at the $\tilde{\zeta}_\mu$ which differ from P only by polynomials from $\text{span } \mathcal{N}$ with coefficients of $O(10^{-6})$.

Assume that the accuracy level of P is better than that. Then we perform one (or more) classical Newton step(s) from each of the $\tilde{\zeta}_\mu$ separately; due to the good starting values, the convergence is extremely rapid. The “exact” zeros of P thus obtained are, rounded to the decimals given:

$$\begin{pmatrix} .0003023 \\ -.0223103 \\ -.56 \cdot 10^{-11} \end{pmatrix}, \begin{pmatrix} .0002978 \\ .0224103 \\ -.52 \cdot 10^{-11} \end{pmatrix}, \begin{pmatrix} .0001285 \\ -.0087233 \\ .0085957 \end{pmatrix}, \begin{pmatrix} .0001216 \\ .0086008 \\ -.0087232 \end{pmatrix}, \begin{pmatrix} .0001305 \\ .0000013 \\ -.0129372 \end{pmatrix}, \begin{pmatrix} .0001362 \\ -.0000013 \\ .0128816 \end{pmatrix}.$$

A comparison shows that our approximate zeros $\tilde{\zeta}_\mu$ have been rather good approximations of the correct zeros ζ_μ of P in its 6-cluster about the origin. Thus our linearization procedure has been successful in a situation with a perturbation level of $O(10^{-4})$ at a 6-fold zero of a system of 3 polynomials in 3 variables.

5 Conclusions

We have shown how one may locate the zeros in an m -cluster of a multivariate polynomial system which arises through a small perturbation of an m -fold zero by an approach which combines *algebraic techniques* and *analytic and numeric reasoning*. We believe that this combination will be necessary and successful also in other applications of algebraic algorithms to problems with not fully accurate data. Our principal tool has been the use of *perturbed Groebner bases* for the representation of ideals near degenerate situations. The motivation for this concept and techniques for its handling have been further explained in a separate publication ([4]).

After a shift of the m -fold zero to the origin, our approach uses only a well-defined set of coefficients of (relatively) low order terms in the polynomials. Therefore, the same approach may immediately be extended to the analysis of a zero cluster about the origin of a system of s *analytic functions* in s variables with a sufficiently large radius of convergence of their Taylor series about the origin. This may be of interest for polynomial systems which contain some trigonometric or exponential terms, even for $s = 1$.

My understanding of the situation has profited a good deal from various discussions with Prof. H.M. Moeller for which I wish to express my sincere gratitude. The complete

procedure for the solution of the *cluster analysis problem* stated at the end of section 1 has been implemented in Maple by my student G. Thallinger in his Diploma Thesis [9]. I wish to thank him for his help in the preparation of this paper, particularly for the running of the example in Section 4 at a time when his program was still in a preliminary state.

6 Appendix: Perturbed Groebner Bases

The following is a *synopsis* of the author's paper [4] on perturbed Groebner bases. It has been included to make the present paper understandable without the knowledge of [4].

For fixed term order, the mapping from a given ideal \mathcal{I} to its (reduced) Groebner basis \mathcal{G} (with leading coefficients 1) is well-defined; \mathcal{G} may be computed from any generating set of \mathcal{I} . In the conceptual context of numerical nonlinear algebra (cf. section 1), the *continuity* of this mapping attains prime importance. For 0-dimensional ideals in \mathbb{P}^s , a *topology* may be introduced in the set of all ideals with a given fixed number m of zeros through the bijective relation between ideals \mathcal{I} and zero sets $V[\mathcal{I}]$ where the natural topology in $\mathbb{C}^{m \cdot s}$ may be used. This topology may be extended to the case of confluent zeros with the aid of the differential conditions at a multiple zero; thus m will always count multiplicities. Informally, continuity of $\mathcal{I} \rightarrow \mathcal{G}$ means that a small change in $V[\mathcal{I}]$ leads to a small change in \mathcal{G} , in particular that it leaves its number of elements and its leading terms unchanged.

But simple examples show that the mapping $V[\mathcal{I}] \rightarrow \mathcal{G}$ has *structural discontinuities*: There exist manifolds RS in $\mathbb{C}^{m \cdot s}$ such that $V_0 \in RS$ implies: In each neighborhood of V_0 there exist V_1, V_2 such that $\mathcal{G}(V_1), \mathcal{G}(V_2)$ have *different normal sets*. This excludes a continuous transition from $\mathcal{G}(V_1)$ to $\mathcal{G}(V_2)$.

The reason for this is immediate: In $\mathbb{C}^s, s > 1$, no set \mathcal{N} of m terms is a *proper interpolation basis* for *all* $V \in (\mathbb{C}^s)^m$. Whichever \mathcal{N} arises as normal set of some \mathcal{I} , the Gram matrix $(\mathcal{N}(V))$ has *rank deficiencies* along certain manifolds $RS \in \mathbb{C}^{m \cdot s}$. As V approaches such a manifold, some coefficients in the associated Groebner basis $\mathcal{G}(V)$ diverge to infinity; on the manifold, a structurally different $\mathcal{G}(V)$, with a different normal set, appears. These manifolds of *representation singularities* are commonly associated with *non-generic* positions of zeros: symmetries, confluences, and the like. Since such degeneracies often occur in applications, this intrinsic shortcoming of classical Groebner bases is deplorable.

However, at least in the *complete intersection case*, our observation which explains the phenomenon also points to a natural way of relief: For *fixed* \mathcal{N} , the Gram matrix $(\mathcal{N}(V))$ depends continuously on V ; therefore, the normal set \mathcal{N}_0 for V_0 on a representation singularity is a proper interpolation basis also in a *full neighborhood* of V_0 . This permits a continuous extension of $\mathcal{G}(V_0)$ into bases for the ideals of all V near V_0 :

Definition: For a fixed term order, consider the Groebner basis $\mathcal{G}_0 = \{g_\kappa, \kappa = 1(1)k\}$ for the ideal \mathcal{I}_0 , with normal set $\mathcal{N}_0 = \{t_\mu, \mu = 1(1)m\}$ and zero set V_0 . A basis $\tilde{\mathcal{G}}$ for an ideal \mathcal{I} with a zero set V close to V_0 is called a *perturbed Groebner basis* of \mathcal{I} if it has the form

$$\tilde{\mathcal{G}} = \{\tilde{g}_\kappa = g_\kappa + \sum_{\mu} \gamma_{\kappa\mu} t_\mu, \quad \kappa = 1(1)k\}. \quad \square \quad (32)$$

The *perturbation vectors* $c_\kappa^T = (\gamma_{\kappa\mu}) \in \mathbb{C}^m$ are defined through *interpolation at V*: If c_κ^T solves (with our previous notation for the normal term vector \mathbf{t})

$$c_\kappa^T \mathbf{t}(\zeta) = g_\kappa(\zeta), \quad \zeta \in V,$$

then $\tilde{\mathcal{G}}(V) = \{g_\kappa - c_\kappa^T \mathbf{t}, \kappa = 1(1)k\}$ is a perturbed Groebner basis of the ideal $\mathcal{I}(V)$ and the c_κ^T depend continuously on V , for V sufficiently close to V_0 .

If V_0 is at a representation singularity, some of the perturbed Groebner bases $\tilde{\mathcal{G}}(V)$ for V near V_0 cannot be genuine Groebner bases; this happens through the introduction into some \tilde{g}_κ of a term t from \mathcal{N} with $t > \text{LT}[g_\kappa]$. However, for sufficiently small c_κ^T , perturbed Groebner bases are nearly as well-suited for computational purposes as classical Groebner bases, – or even more so, considering their continuous behavior.

For ideals with m zeros, a neighborhood of V_0 is a region in $(\mathbb{C}^s)^m$ while $\tilde{\mathcal{G}}$ of (32) contains k coefficient vectors $c_\kappa^T \in \mathbb{C}^m$. Thus, for $k = s$, the neighborhood of V_0 is mapped onto a full neighborhood of the origin in the coefficient space. For $k > s$, on the other hand, the c_κ^T which may appear in $\tilde{\mathcal{G}}$ must lie in an $m s$ -dimensional manifold within the $m k$ -dimensional neighborhood of the origin. This manifold SZ is defined by the extension of the *syzygies* of \mathcal{G}_0 which pose restrictions on the values which the c_κ^T may take.

Our use of this bijective mapping between a neighborhood of V_0 and a neighborhood of the origin (on SZ) in the coefficient space is explained by the following diagram:

$$\begin{array}{ccccccc} \text{system } P_0 & \longrightarrow & \text{zero set } V_0 & \longleftrightarrow & \text{mult. tables} & \longleftrightarrow & \text{Groebner basis } \mathcal{G}_0 \\ \text{system } P & \longrightarrow & \text{zero set } V & \longleftrightarrow & \text{mult. tables } A^* & \longleftrightarrow & \text{pert. Groebner basis } \mathcal{G}^*(V) \\ & \searrow & \text{appr. zero set } \tilde{V} & \longleftarrow & \text{appr. mult. tables } \tilde{A} & \longleftarrow & \text{appr. pert. Gr. basis } \tilde{\mathcal{G}}_{\tilde{c}} \\ & & \searrow & & \text{linearization of } P \rightarrow \mathcal{G}^*(V) \text{ yields } \tilde{c}_\kappa & & \nearrow \end{array}$$

In order to obtain V from P , we should find $\mathcal{G}^*(V)$ directly from P and then obtain V via the associated multiplication tables. The computation of the perturbation coefficients in $\mathcal{G}^*(V)$ is generally too involved; since we expect them to be small for P close to P_0 we *linearize* the mapping from P to $\mathcal{G}^*(V)$ about $c_\kappa^T = 0$ and proceed with the approximate perturbation vectors \tilde{c}_κ^T thus obtained towards an approximation \tilde{V} of V . In the computation of the \tilde{c}_κ^T , we restrict them to the $m \times s$ -dimensional tangent space of the manifold SZ at the origin which is obtained by linearization of the syzygies of $\mathcal{G}^*(V)$ w.r.t. the c_κ^T at 0.

The computational details have been explained in section 3.3; cf. also section 4. A more concise technical description of the preceeding explanations is to be found in [4].

References

- [1] W. Kahan: Conserving Confluence Curbs Ill-condition; Comp. Science UC Berkeley Tech.Rep. 6, 1972
- [2] V. Hribernik: Sensitivity of Algebraic Algorithms, PhD Thesis, TU Vienna, 1995

- [3] V. Hribernig, H.J. Stetter: Detection and Validation of Clusters of Polynomial Zeros, submitted to J.Symb.Comp.
- [4] H.J. Stetter: Perturbed Groebner Basis Representations, in preparation
- [5] H.M. Moeller, H.J. Stetter: Multivariate Polynomial Equations with Multiple Zeros Solved by Matrix Eigenproblems, Numer. Math. **70** (1995) 311–329
- [6] M.G. Marinari, H.M. Moeller, T. Mora: Groebner Bases of Ideals Given by Dual Bases, Proceedings of ISSAC 91, 55–63
- [7] M.G. Marinari, H.M. Moeller, T. Mora: Groebner Duality and Multiplicities in Polynomial System Solving, to appear in
- [8] E. Anderson et.al.: LAPACK Users’ Guide, 2nd Ed., SIAM Publ., 1995
- [9] G. Thallinger: Analysis of Zero Clusters in Multivariate Polynomial Systems, Diploma Thesis, TU Vienna, 1996
- [10] Th. Becker, V. Weispfenning: Groebner Bases: A Computational Approach to Commutative Algebra, Springer Grad. Texts in Math. **141**, 1993
- [11] H.J. Stetter: Multivariate Polynomial Equations as Matrix Eigenproblems, WSSIA **2** (1993) 355–371

Part of the research reported in this paper was done while the author spent a sabbatical term at the Computer Science Department of Cornell University.