

EXISTENCE VERIFICATION FOR SINGULAR ZEROS OF NONLINEAR SYSTEMS*

R. BAKER KEARFOTT[†], JIANWEI DIAN[‡], AND A. NEUMAIER[§]

Abstract. Computational fixed point theorems can be used to automatically verify existence and uniqueness of a solution to a nonlinear system of equations $F(x) = 0$, $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$ within a given region \mathfrak{x} of n -space. But such computations succeed only when the Jacobi matrix $F'(x)$ is nonsingular everywhere in \mathfrak{x} . However, in many practical problems, the Jacobi matrix is singular, or nearly so, at the solution x^* , $F(x^*) = 0$. In such cases, arbitrarily small perturbations of the problem result in problems $\tilde{F}(x) = 0$ either with *no* solutions in \mathfrak{x} or with *more than one* solution in \mathfrak{x} ; thus no general computational technique can prove existence and uniqueness. This leads to a fundamental philosophical problem: “What is meant by existence and uniqueness in such cases?”

Here, an interpretation of verification is given in the singular context: proof that a given number of true solutions exist within a region in complex space containing \mathfrak{x} .

Proof that a given number of true solutions exist within a given region of complex space is possible by computation of the topological degree, but such computations heretofore have required a global search on the $(n - 1)$ -dimensional boundary of an n -dimensional region. Here, it is observed that preconditioning leads to a system of equations whose topological degree can be computed with a much lower-dimensional search. Formulas are given for this computation, and the special case of rank-defect one is studied, both theoretically and empirically.

Key words. nonlinear systems, interval computations, verified computations, singularities, topological degree

AMS subject classifications. 65G10, 65H10

1. Introduction. Given an approximate solution \tilde{x} to a nonlinear system of equations $F(x) = 0$, $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$, it is useful in various contexts to construct bounds around \tilde{x} in which it is proven that there exists a unique solution x^* , $F(x^*) = 0$. For continuously differentiable F for which the Jacobian $\det(F'(x^*)) \neq 0$ and for which that Jacobian is well conditioned, interval computations have no trouble proving that there is a unique solution within small boxes with x^* reasonably near the center; see [5, 9, 10]. However, if $F'(x^*)$ is ill-conditioned or singular, such computations necessarily must fail. In this singular situation, what is meant by “unique solution” can be defined in several ways. A particular interpretation may be most appropriate in a particular practical context. Also, given a particular interpretation, various computational procedures can be developed to verify existence and uniqueness. In this paper, one interpretation and its corresponding computational procedure are considered.

1.1. Notation. We assume familiarity with the fundamentals of interval arithmetic; see [9, 10] for an introduction in the present context. (The works [2, 5, 11] also contain introductory material.)

Throughout, scalars and vectors will be denoted by lower case, while matrices will be denoted by upper case. Intervals, interval vectors (also called “boxes”) and interval matrices will be denoted by boldface. For instance, $\mathbf{x} = (x_1, \dots, x_n)$ denotes

*This work was supported by National Science Foundation grant DMS-9701540.

[†]Department of Mathematics, University of Louisiana at Lafayette, Lafayette, LA 70504, USA (rbk@usl.edu).

[‡]Department of Mathematics, University of Louisiana at Lafayette, Lafayette, LA 70504, USA (dian@usl.edu).

[§]Institut für Mathematik, Universität Wien, Strudhofgasse 4, A-1050 Wien, Austria (neum@cma.univie.ac.at).

an interval vector, $A = (a_{i,j})$ denotes a point matrix, and $\mathbf{A} = (\mathbf{a}_{i,j})$ denotes an interval matrix. Real n -space will be denoted by \mathbb{R}^n , while the set of n -dimensional interval vectors, i.e. n -dimensional boxes, will be denoted by \mathbb{IR}^n . Similarly, complex n -space will be denoted by \mathbb{C}^n , while the space of n -dimensional complex interval vectors will be denoted by \mathbb{IC}^n . The midpoint of an interval or interval vector \mathbf{x} will be denoted by $m(\mathbf{x})$. The non-oriented boundary of a box \mathbf{x} will be denoted by $\partial\mathbf{x}$ while its oriented boundary will be denoted by $b(\mathbf{x})$. (See §2.)

1.2. Traditional Computational Existence and Uniqueness. Computational existence and uniqueness verification rests on interval versions of Newton's method. Typically, such computations can be described as evaluation of a related interval operator $\mathbf{G}(\mathbf{x})$; $\mathbf{G}(\mathbf{x}) \subseteq \mathbf{x}$ then implies existence and uniqueness of $F(x) = 0$ within \mathbf{x} . To describe these, we review

DEFINITION 1.1. ([10, p. 174], etc.) *Let $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$. The matrix \mathbf{A} is said to be a Lipschitz matrix for F over \mathbf{x} provided, for every $x \in \mathbf{x}$ and $y \in \mathbf{x}$, $F(x) - F(y) = A(x - y)$ for some $A \in \mathbf{A}$.*

Most interval Newton methods for $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$, abstractly, are of the general form

$$(1.1) \quad \tilde{\mathbf{x}} = \mathbf{N}(F; \mathbf{x}, \tilde{x}) = \tilde{x} + \mathbf{v},$$

where \mathbf{v} is computed to contain the solution set to the interval linear system

$$(1.2) \quad \mathbf{A}\mathbf{v} = -F(\tilde{x}),$$

and where, for initial uniqueness verification, \mathbf{A} is generally a Lipschitz matrix¹ for F over the box (interval vector) \mathbf{x} and $\tilde{x} \in \mathbf{x}$ is a guess point. We sometimes write $\mathbf{F}'(\mathbf{x})$ in place of \mathbf{A} , since the matrix can be an interval extension of the Jacobi matrix of F . Uniqueness verification traditionally depends on regularity of the matrix \mathbf{A} . We have

LEMMA 1.2. ([9, 10]) *Suppose $\tilde{\mathbf{x}} = \tilde{x} + \mathbf{v}$ is the image under the interval Newton method (formula (1.1)), where \mathbf{v} is computed by any method that bounds the solution set to the interval linear system (1.2), and $\tilde{\mathbf{x}} \subseteq \mathbf{x}$. Then \mathbf{A} is regular.*

The method of bounding the solution set of equation (1.2) to be considered here is the interval Gauss–Seidel method, defined by:

DEFINITION 1.3. *The preconditioned interval Gauss–Seidel image $\mathbf{GS}(F; \mathbf{x}, \tilde{x})$ of a box \mathbf{x} is defined as $\mathbf{GS}(F; \mathbf{x}, \tilde{x}) \equiv (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n)$, where $\tilde{\mathbf{x}}_i$ is defined sequentially for $i = 1$ to n by:*

$$\tilde{\mathbf{x}}_i \equiv \mathbf{x}_i \cap \left(\tilde{x}_i - \mathbf{N}_i / (Y_i \mathbf{A}_i) \right),$$

where

$$\mathbf{N}_i = Y_i F(\tilde{x}) + \sum_{j=1}^{i-1} Y_i \mathbf{A}_j (\tilde{\mathbf{x}}_j - \tilde{x}_j) + \sum_{j=i+1}^n Y_i \mathbf{A}_j (\mathbf{x}_j - \tilde{x}_j),$$

and where $\tilde{x} = (\tilde{x}_1, \dots, \tilde{x}_n)^T$ is an initial guess point, $Y\mathbf{A} \in \mathbb{IR}^{n \times n}$ and $YF(\tilde{x})$ are the matrix and right-hand-side vector for the preconditioned interval system $Y\mathbf{A}(x - \tilde{x}) = -YF(\tilde{x})$, $Y \in \mathbb{R}^{n \times n}$ is a point preconditioning matrix, Y_i denotes the i -th row of Y , and \mathbf{A}_j denotes the j -th column of \mathbf{A} .

¹However, see [9, 12] for techniques for using slope matrices.

Lemma 1.2 applies when $\mathbf{N}(F; \mathbf{x}, \tilde{x}) = \mathbf{GS}(F; \mathbf{x}, \tilde{x})$, provided we specify that $\mathbf{GS}(F; \mathbf{x}, \tilde{x})$ be in the interior $\text{int}(\mathbf{x})$ of \mathbf{x} . (We must specify the interior because of the intersection step in Definition 1.3.) In particular, we have

THEOREM 1.4. ([9, 10]) *Suppose $F : \mathbf{x} \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ and \mathbf{A} is a Lipschitz matrix such as an interval extension $\mathbf{F}'(\mathbf{x})$ of the Jacobi matrix. If \tilde{x} is the image under an interval Newton method as in formula (1.1) and $\tilde{x} \subset \text{int}(\mathbf{x})$, then there is a unique $x^* \in \mathbf{x}$ with $F(x^*) = 0$. Various authors have proven Theorem 1.4; see [9, 10]. In particular, Miranda's theorem can be used to easily prove Theorem 1.4 for $\mathbf{N}(F; \mathbf{x}, \tilde{x}) = \mathbf{GS}(F; \mathbf{x}, \tilde{x})$; see [9, p. 60].*

EXAMPLE 1. *Take*

$$\begin{aligned} f_1(x_1, x_2) &= x_1^2 - x_2, \\ f_2(x_1, x_2) &= x_1 - x_2^2, \end{aligned}$$

and $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)^T = ([-0.1, 0.1], [-0.1, 0.3])^T$. There is a unique root $x^* = (0, 0)^T$ of $F = (f_1, f_2)^T$ within \mathbf{x} . In Example 1, if $\tilde{x} = (0, .1)^T$, then $F(\tilde{x}) = (-0.1, -0.01)^T$, and an interval extension of the Jacobi matrix is

$$\mathbf{F}'(\mathbf{x}) = \begin{pmatrix} 2\mathbf{x}_1 & -1 \\ 1 & -2\mathbf{x}_2 \end{pmatrix} = \begin{pmatrix} [-.2, .2] & -1 \\ 1 & [-0.6, 0.2] \end{pmatrix}.$$

If the preconditioning matrix Y is taken to be the inverse of the midpoint matrix of $\mathbf{F}'(\mathbf{x})$, then

$$Y = \{\text{m}(\mathbf{F}'(\mathbf{X}))\}^{-1} = \begin{pmatrix} -0.2 & 1 \\ -1 & 0 \end{pmatrix}.$$

We obtain, rounded out to four digits,

$$\mathbf{GS}(F; \mathbf{x}, \tilde{x}) = ([-0.07292, 0.09375], [-0.01875, 0.01875])^T \subset \text{int}(\mathbf{x}).$$

Therefore, this computation proves that there is a unique solution to $F(x) = 0$ within \mathbf{x} .

Inclusion in the interval Gauss–Seidel method is made possible because the inverse midpoint preconditioner reduces the interval Jacobi matrix to approximately a diagonal matrix.

1.3. Singularities: Philosophical Considerations. Theorem 1.4 is applicable only when the matrix \mathbf{A} is regular, i.e. when \mathbf{A} does not contain singular matrices. But consider

EXAMPLE 2. *Take*

$$\begin{aligned} f_1(x_1, x_2) &= x_1^2 - x_2, \\ f_2(x_1, x_2) &= x_1^2 + x_2, \end{aligned}$$

and $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)^T = ([-0.001, 0.001], [-0.001, 0.001])^T$. Even though there is a unique root $x^* = (0, 0)^T$ of $F = (f_1, f_2)^T$ within \mathbf{x} when F is as in Example 2, the interval Gauss–Seidel method cannot prove this, since the Jacobi matrix $F'(x^*)$ is singular. In fact, the interval Jacobi matrix is computed to be

$$\mathbf{F}'(\mathbf{x}) = \begin{pmatrix} 2\mathbf{x}_1 & -1 \\ 2\mathbf{x}_1 & 1 \end{pmatrix} = \begin{pmatrix} [-0.002, 0.002] & -1 \\ [-0.002, 0.002] & 1 \end{pmatrix},$$

and the midpoint

$$m(\mathbf{F}'(\mathbf{x})) = \begin{pmatrix} 0 & -1 \\ 0 & 1 \end{pmatrix},$$

matrix the inverse of which is often used as the preconditioner matrix Y , is not invertible².

Symbolic methods can be used to show that Example 2 has a unique solution in a small region containing $x_1 = 0$, $x_2 = 0$. However, arbitrarily small perturbations of the problem result in either no solutions or two solutions. Consider

EXAMPLE 3. *Take*

$$\begin{aligned} f_1(x_1, x_2) &= x_1^2 - x_2, \\ f_2(x_1, x_2) &= x_1^2 + x_2 + \epsilon, \end{aligned}$$

and $\mathbf{x} = (x_1, x_2)^T = ([-0.001, 0.001], [-0.001, 0.001])^T$. Here, $|\epsilon|$ is very small.

The system in Example 3 has two solutions for $\epsilon < 0$ and no solutions for $\epsilon > 0$. But roundout in computer arithmetic and, perhaps, uncertainties in the system itself due to modelling or measurement uncertainties, make it impossible to distinguish systems such as in Example 3 for different ϵ , especially when computer arithmetic is used as part of the verification process. In such instances, a new interpretation needs to be given to existence / uniqueness. One expects the appropriate interpretation to be application-dependent in general, but the following is possible: **verify that the system has an exact number of solutions within a larger space containing the original space**. This paper will be based on this interpretation.

To be specific, we can extend an n -dimensional box in \mathbb{R}^n to an n -dimensional box in \mathbb{C}^n by adding a small imaginary part to each variable. If the system can be extended to an analytic function in complex n -space (or if it can be extended to a function that can be approximated by an analytic function), then the *topological degree* gives the number of solutions, counting multiplicities, within the small region in complex space. (See §2 for an explanation of multiplicity.) When the imaginary parts are small, we can view the complex roots as approximating a multiple real root. For example, the degree of the system in Example 3 within an extended box in complex space can be computed to be 2, regardless of whether ϵ is negative, positive or zero. (See the numerical results in §8.) The topological degree corresponds roughly to algebraic degree in one dimension; for example, the degree of z^n in a small region in \mathbb{C}^1 containing 0 is n .

1.4. Organization of This Paper. A review of properties of the topological degree, to be used later, appears in §2. The issue of preconditioning appears in §3. Construction of the box in the complex space appears in §4.

Several algorithms have previously been proposed for computing the topological degree [1, 7, 14], but these require computational effort equivalent to finding all solutions to $4n$ ($2n-1$)-dimensional nonlinear systems within a given box, or worse. In §5, a reduction is proposed that allows computation of the topological degree with a search in a space of dimension equal to the rank defect of the Jacobian matrix. A theorem is proven that further simplifies the search.

²Alternate preconditioners can nonetheless be computed; see [9]. However, it can be shown that uniqueness cannot be proven in this case; see [9, 10].

In §6, the actual algorithm is presented and its computational complexity is given. Test problems and the test environment are described in §7. Numerical results appear in §8. Future directions appear in §9.

2. Review of Some Elements of Degree Theory. The topological degree or Brouwer degree, well-known within algebraic topology and nonlinear functional analysis, is both a generalization of the concept of a sign change of a one-dimensional continuous function and of the winding number for analytic functions. It can be used to generalize the concept of multiplicity of a root. The fundamentals will not be reviewed here, but we refer to [3, 4, 7]. We only present the material we need.

Here, we explain what we mean by “multiplicity”. Actually, there is a more general concept *index* (see [4, Chapter I]) for an isolated zero. The topological degree is equal to the sum of the indices of zeros in the domain. The index is always positive in our context. So, we use multiplicity as an alternative term for index since it’s more suggestive.

Suppose that $F : \mathbf{D} \subset \mathbb{C}^n \rightarrow \mathbb{C}^n$ is analytic. Then the real and imaginary components of F and its argument $z \in \mathbb{C}^n$ may be viewed as real components in \mathbb{R}^{2n} . Let $z = x + iy$ and $F(z) = u(x, y) + iv(x, y)$, where $x = (x_1, \dots, x_n), y = (y_1, \dots, y_n), u(x, y) = (u_1(x, y), \dots, u_n(x, y))$ and $v(x, y) = (v_1(x, y), \dots, v_n(x, y))$. We can define $\tilde{\mathbf{D}}$ by

$$\tilde{\mathbf{D}} \equiv \{(x_1, y_1, \dots, x_n, y_n) | (x_1 + iy_1, \dots, x_n + iy_n) \in \mathbf{D}\}$$

and $\tilde{F} : \tilde{\mathbf{D}} \subset \mathbb{R}^{2n} \rightarrow \mathbb{R}^{2n}$ by

$$\tilde{F} = (u_1, v_1, \dots, u_n, v_n).$$

Then, we have the following property of topological degree $d(\tilde{F}, \tilde{\mathbf{D}}, 0)$, and relationships between $d(\tilde{F}, \tilde{\mathbf{D}}, 0)$ and the solutions of the system $F(z) = 0$ in \mathbf{D} .

THEOREM 2.1. *Suppose $F : \mathbf{D} \subset \mathbb{C}^n \rightarrow \mathbb{C}^n$ is analytic, with $F(z) \neq 0$ for any $z \in \partial\mathbf{D}$, and suppose $\tilde{\mathbf{D}}$ and $\tilde{F} : \tilde{\mathbf{D}} \rightarrow \mathbb{R}^{2n}$ are defined as above. Then*

1. $d(\tilde{F}, \tilde{\mathbf{D}}, 0) \geq 0$.
2. $d(\tilde{F}, \tilde{\mathbf{D}}, 0) > 0$ if and only if there is a solution $z^* \in \mathbf{D}$, $F(z^*) = 0$.
3. $d(\tilde{F}, \tilde{\mathbf{D}}, 0)$ is equal to the number of solutions $z^* \in \mathbf{D}$, $F(z^*) = 0$, counting multiplicities.
4. If the Jacobi matrix $F'(z^*)$ is non-singular at every $z^* \in \mathbf{D}$ with $F(z^*) = 0$, then $d(\tilde{F}, \tilde{\mathbf{D}}, 0)$ is equal to the number of solutions $z^* \in \mathbf{D}$, $F(z^*) = 0$.

To obtain Theorem 2.1, we need to notice that $F(z^*) = 0$ for $z^* \in \mathbf{D}$ is equivalent to $\tilde{F}(x^*, y^*) = 0$ for $(x^*, y^*) \in \tilde{\mathbf{D}}$, where $z^* = x^* + iy^*$. In Theorem 2.1, the first conclusion is easily obtained by knowing $|d(F_R, \mathbf{D}_R, 0)| \leq d(\tilde{F}, \tilde{\mathbf{D}}, 0)$, where F_R is function $F(z)$ with z in a corresponding n -dimensional real domain $\mathbf{D}_R \in \mathbb{R}^n$. (See [4, p. 49].) The third conclusion is actually Theorem (5.2) in Chapter I of [4], considering the definition of topological degree. The fourth conclusion can be obtained by the third conclusion, Theorem (7.2) and Lemma (9.3-2) in Chapter I of [4]. As for the second conclusion, it’s easy to see that $d(\tilde{F}, \tilde{\mathbf{D}}, 0) > 0$ implies there is a solution $(x^*, y^*) \in \tilde{\mathbf{D}}$, $\tilde{F}(x^*, y^*) = 0$. (See Existence Theorem (6.6) in Chapter I of [4].) The reverse is obtained by noticing the third conclusion and the fact that the multiplicities are positive.

The following three theorems will lead to the degree computation formula in Theorem 5.1 in §5, the formula used in our computational scheme.

THEOREM 2.2. *(See [13, §4.2].) Let \mathbf{D} be an n -dimensional connected, oriented region in \mathbb{R}^n and $F = (f_1, \dots, f_n)$, where $f_k, k = 1, \dots, n$, are continuous functions*

defined in \mathbf{D} . Assume $F \neq 0$ on the oriented boundary $b(\mathbf{D})$ of \mathbf{D} , $b(\mathbf{D})$ can be subdivided into a finite number of closed, connected $(n-1)$ -dimensional oriented subsets β_{n-1}^k , $k = 1, \dots, r$, and there is p , $1 \leq p \leq n$, such that:

1. $F_{-p} \equiv (f_1, \dots, f_{p-1}, f_{p+1}, \dots, f_n) \neq 0$ on the oriented boundary $b(\beta_{n-1}^k)$ of β_{n-1}^k , $k = 1, \dots, r$; and
2. f_p has the same sign at all solutions of $F_{-p} = 0$, if there are any, on β_{n-1}^k , $1 \leq k \leq r$.

Choose $s \in \{-1, +1\}$ and let $K_0(s)$ denote the subset of the integers $k \in \{1, \dots, r\}$ such that $F_{-p} = 0$ has solutions on β_{n-1}^k and $\text{sgn}(f_p) = s$ at each of those solutions. Then

$$d(F, \mathbf{D}, 0) = (-1)^{p-1} s \sum_{k \in K_0(s)} d(F_{-p}, \beta_{n-1}^k, 0).$$

The formula in Theorem 2.2 is a combination of formulas (4.15) and (4.16) in [13]. The orientation of \mathbf{D} is positive and the orientations of β_{n-1}^k , whether positive or negative, are induced by the orientation of \mathbf{D} . If we assume that the Jacobi matrices of F_{-p} are non-singular at all solutions of $F_{-p} = 0$ on β_{n-1}^k , then

$$d(F_{-p}, \beta_{n-1}^k, 0) = t(\beta_{n-1}^k) \sum_{\substack{x \in \beta_{n-1}^k \\ F_{-p}(x)=0}} \text{sgn}(JF_{-p}(x)),$$

where $t(\beta_{n-1}^k) = +1$ or -1 depending on whether β_{n-1}^k has positive orientation or negative orientation, and $JF_{-p}(x)$ is the determinant of the Jacobi matrix of F_{-p} at x . (See Theorem (5.2) and (7.2) in Chapter I of [4].) So, under this assumption, we can simplify the formula in Theorem 2.2.

THEOREM 2.3. *Suppose that all the conditions of Theorem 2.2 are satisfied and the Jacobi matrices of F_{-p} are non-singular at all solutions of $F_{-p} = 0$ on β_{n-1}^k for each $k \in K_0(s)$. Then*

$$d(F, \mathbf{D}, 0) = (-1)^{p-1} s \sum_{k \in K_0(s)} t(\beta_{n-1}^k) \sum_{\substack{x \in \beta_{n-1}^k \\ F_{-p}(x)=0}} \text{sgn}(JF_{-p}(x)),$$

where $t(\beta_{n-1}^k) = +1$ or -1 depending on whether β_{n-1}^k has positive orientation or negative orientation, and $JF_{-p}(x)$ is the determinant of the Jacobi matrix of F_{-p} at x .

In our context, the region \mathbf{D} is an n -dimensional box

$$\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n),$$

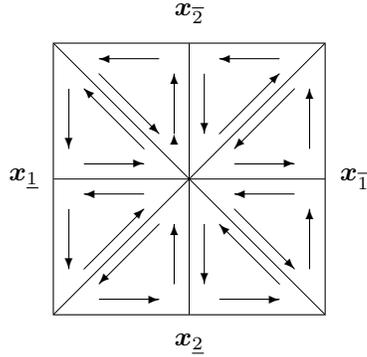
where $n \geq 2$ and $\mathbf{x}_k = [\underline{x}_k, \bar{x}_k]$. The boundary $\partial \mathbf{x}$ of \mathbf{x} consists of $2n(n-1)$ -dimensional boxes

$$\mathbf{x}_k \equiv (\mathbf{x}_1, \dots, \mathbf{x}_{k-1}, \underline{x}_k, \mathbf{x}_{k+1}, \dots, \mathbf{x}_n)$$

and

$$\mathbf{x}_k^- \equiv (\mathbf{x}_1, \dots, \mathbf{x}_{k-1}, \bar{x}_k, \mathbf{x}_{k+1}, \dots, \mathbf{x}_n),$$

where $k = 1, \dots, n$.

FIG. 2.1. Positive center decomposition of \mathbf{x} when $n = 2$.

Before introducing the next theorem, for completeness in the exposition, we construct the so-called positive or negative *center decomposition* of \mathbf{x} (although such decompositions are undoubtedly known among experts). This center decomposition shows us how to specify orientations of the faces of \mathbf{x} . We follow the convention that an n -simplex $\langle P_0, P_1, \dots, P_n \rangle$ in \mathbb{R}^n has positive or negative orientation if the determinant $\Delta_n((1, P_0), (1, P_1), \dots, (1, P_n))$ is positive or negative. Here, $P_i = (x_{i1}, \dots, x_{in}), i = 0, 1, \dots, n$, are $n + 1$ points in \mathbb{R}^n represented by row vectors, and $(1, P_i) = (1, x_{i1}, \dots, x_{in}), i = 0, 1, \dots, n$, are $(n + 1)$ -dimensional vectors that form the $n + 1$ rows of the determinant. See [3, 4, 7] for further details.

We construct the center decompositions of \mathbf{x} inductively. Without loss of generality, assume $\mathbf{x}_k = [-1, +1], k = 1, \dots, n$. So, the center of \mathbf{x} is $(0, \dots, 0)$. When $n = 2$ the positive center decomposition is shown in Figure 2.1, in which $P_0 = (0, 0)$. It's easy to check that each of the $2^2 \cdot 2!$ 2-simplexes has positive orientation. So, \mathbf{x} has positive orientation. Sometimes we say the orientation is $+1$ or -1 if it's positive or negative. It's clear that the orientations of $\mathbf{x}_1, \mathbf{x}_1, \mathbf{x}_2$ and \mathbf{x}_2 are $-1, +1, +1$ and -1 . The negative center decomposition of \mathbf{x} can be obtained by taking all the simplexes with the negative orientation. Sometimes we say a center decomposition with sign $+1$ or -1 if it's positive or negative, respectively.

Suppose we already have center decompositions of $(n - 1)$ -dimensional ($n > 2$) boxes with each decomposition consisting of $2^{n-1} \cdot (n - 1)!$ $(n - 1)$ -simplexes. We construct the positive center decomposition of an n -dimensional box \mathbf{x} as follows. We take the center decomposition with sign $(-1)^k$ of the $(n - 1)$ -dimensional box \mathbf{x}_k and we take the center decomposition with sign $(-1)^{k+1}$ of the $(n - 1)$ -dimensional box \mathbf{x}_k , where $k = 1, \dots, n$.

We use

$$\langle Q_0^{(k,-1)}, Q_1^{(k,-1)}, \dots, Q_{n-1}^{(k,-1)} \rangle$$

or

$$\langle Q_0^{(k,+1)}, Q_1^{(k,+1)}, \dots, Q_{n-1}^{(k,+1)} \rangle$$

to denote an $(n - 1)$ -simplex in the center decomposition of \mathbf{x}_k or \mathbf{x}_k , where

$$Q_i = (x_{i1}, \dots, x_{i(k-1)}, x_{i(k+1)}, \dots, x_{in})$$

and

$$Q_i^{(k,\pm 1)} = (x_{i1}, \dots, x_{i(k-1)}, \pm 1, x_{i(k+1)}, \dots, x_{in}).$$

So,

$$\begin{aligned} & \text{sgn}(\Delta_{n-1}((1, Q_0), (1, Q_1), \dots, (1, Q_{n-1}))) \\ &= (-1)^k \text{ on } \mathbf{x}_{\underline{k}} \text{ and } (-1)^{k+1} \text{ on } \mathbf{x}_{\overline{k}}. \end{aligned}$$

For each $(n-1)$ -simplex

$$\langle Q_0^{(k,-1)}, Q_1^{(k,-1)}, \dots, Q_{n-1}^{(k,-1)} \rangle$$

or

$$\langle Q_0^{(k,+1)}, Q_1^{(k,+1)}, \dots, Q_{n-1}^{(k,+1)} \rangle,$$

we add the center $P_0 = (0, 0, \dots, 0)$ of \mathbf{x} as the first point to form an n -simplex

$$\langle P_0, Q_0^{(k,-1)}, Q_1^{(k,-1)}, \dots, Q_{n-1}^{(k,-1)} \rangle$$

or

$$\langle P_0, Q_0^{(k,+1)}, Q_1^{(k,+1)}, \dots, Q_{n-1}^{(k,+1)} \rangle.$$

It's easy to see that there are altogether $2^n \cdot n!$ n -simplexes. Next we verify that each of those n -simplexes has the positive orientation.

$$\begin{aligned} & \Delta_n((1, P_0), (1, Q_0^{(k,-1)}), (1, Q_1^{(k,-1)}), \dots, (1, Q_{n-1}^{(k,-1)})) \\ &= \Delta_{n-1}(Q_0^{(k,-1)}, Q_1^{(k,-1)}, \dots, Q_{n-1}^{(k,-1)}) \\ &= (-1)^k \Delta_{n-1}((1, Q_0), (1, Q_1), \dots, (1, Q_{n-1})). \\ & \Delta_n((1, P_0), (1, Q_0^{(k,+1)}), (1, Q_1^{(k,+1)}), \dots, (1, Q_{n-1}^{(k,+1)})) \\ &= \Delta_{n-1}(Q_0^{(k,+1)}, Q_1^{(k,+1)}, \dots, Q_{n-1}^{(k,+1)}) \\ &= (-1)^{k-1} \Delta_{n-1}((1, Q_0), (1, Q_1), \dots, (1, Q_{n-1})). \end{aligned}$$

So,

$$\begin{aligned} & \text{sgn}(\Delta_n((1, P_0), (1, Q_0^{(k,-1)}), (1, Q_1^{(k,-1)}), \dots, (1, Q_{n-1}^{(k,-1)}))) \\ &= (-1)^k (-1)^k = +1, \text{ and} \\ & \text{sgn}(\Delta_n((1, P_0), (1, Q_0^{(k,+1)}), (1, Q_1^{(k,+1)}), \dots, (1, Q_{n-1}^{(k,+1)}))) \\ &= (-1)^{k-1} (-1)^{k+1} = +1. \end{aligned}$$

This means we actually do get a positive center decomposition of \mathbf{x} . We can get a negative center decomposition by taking all the $2^n \cdot n!$ n -simplexes with negative orientation.

From the above we know that if we decompose \mathbf{x} by the positive center decomposition, then \mathbf{x} has positive orientation, the induced orientation of $\mathbf{x}_{\underline{k}}$ is $(-1)^k$, and the induced orientation of $\mathbf{x}_{\overline{k}}$ is $(-1)^{k+1}$. $b(\mathbf{x})$ can be divided into $\mathbf{x}_{\underline{k}}$ and $\mathbf{x}_{\overline{k}}$, $k = 1, \dots, n$, with the associated orientations. So, $F \neq 0$ on $b(\mathbf{x})$ is the same as $F \neq 0$ on $\partial \mathbf{x}$. Similarly, $F_{-p}(x) = 0$ on $b(\mathbf{x}_{\underline{k}})$ or $b(\mathbf{x}_{\overline{k}})$ is the same as $F_{-p}(x) = 0$ on $\partial \mathbf{x}_{\underline{k}}$ or $\partial \mathbf{x}_{\overline{k}}$. Let $\underline{K}_0(s)$ denote the subset of the integers $k \in \{1, \dots, n\}$ such that $F_{-p} = 0$ has solutions on $\mathbf{x}_{\underline{k}}$ and $\text{sgn}(f_p) = s$ at all the solutions, and $\overline{K}_0(s)$ denote the subset of the integers $k \in \{1, \dots, n\}$ such that $F_{-p} = 0$ has solutions on $\mathbf{x}_{\overline{k}}$ and $\text{sgn}(f_p) = s$ at all the solutions, where $s \in \{-1, +1\}$. Then, by Theorem 2.3, we have

THEOREM 2.4. *Suppose $F \neq 0$ on $\partial \mathbf{x}$, and there is p , $1 \leq p \leq n$, such that:*

1. $F_{-p} \neq 0$ on $\partial \mathbf{x}_{\underline{k}}$ or $\partial \mathbf{x}_{\overline{k}}$, $k = 1, \dots, n$;

$$Y\mathbf{F}'(\mathbf{x}) = \begin{pmatrix} 1 & 0 & \dots & 0 & \overbrace{* \dots *}^p \\ 0 & 1 & 0 \dots & 0 & * \dots * \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & 0 & 1 & * \dots * \\ 0 & \dots & 0 & 0 & 0 \dots 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 0 & 0 \dots 0 \end{pmatrix}.$$

FIG. 3.1. A singular system of rank $n - p$ preconditioned with an incomplete LU factorization, where “*” represents a non-zero element.

2. f_p has the same sign at all solutions of $F_{-p} = 0$, if there are any, on \mathbf{x}_k or \mathbf{x}_k^- , $1 \leq k \leq n$; and
3. the Jacobi matrices of F_{-p} are non-singular at all solutions of $F_{-p} = 0$ on $\partial\mathbf{x}$.

Then

$$\begin{aligned} & d(F, \mathbf{x}, 0) \\ &= (-1)^{p-1} s \left\{ \sum_{k \in \underline{K_0}(s)} (-1)^k \sum_{\substack{x \in \mathbf{x}_k \\ F_{-p}(x)=0}} \operatorname{sgn} \left| \frac{\partial F_{-p}}{\partial x_1 x_2 \dots x_{k-1} x_{k+1} \dots x_n} (x) \right| \right. \\ & \quad \left. + \sum_{k \in \overline{K_0}(s)} (-1)^{k+1} \sum_{\substack{x \in \mathbf{x}_k^- \\ F_{-p}(x)=0}} \operatorname{sgn} \left| \frac{\partial F_{-p}}{\partial x_1 x_2 \dots x_{k-1} x_{k+1} \dots x_n} (x) \right| \right\}. \end{aligned}$$

3. On Preconditioning. The inverse midpoint preconditioner approximately *diagonalizes* the interval Jacobi matrix, when $F'(x^*)$ is non-singular (and well-enough conditioned). This preconditioner can be computed with Gaussian elimination with partial pivoting. We can compute (to within a series of row permutations) an LU factorization of the midpoint matrix $m(\mathbf{F}'(\mathbf{X}))$. The factors L and U may then be applied to actually precondition the interval linear system.

When the rank of $F'(x^*)$ is $n - p$ for some $p > 0$, then Gaussian elimination with full pivoting can be used to reduce $\mathbf{F}'(\mathbf{x})$ to approximately the pattern shown in Figure 3.1. Actually, an incomplete factorization based on full pivoting will put the system into a pattern that resembles a permutation of the columns of the pattern in Figure 3.1. However, for notational simplicity, there is no loss here in assuming exactly the form in Figure 3.1.

In the analysis to follow, we assume that the system has already been preconditioned, so that it is, to within second-order terms with respect to $w(\mathbf{x})$, of the form in Figure 3.1. In the rest of this paper, we concentrate on the case $p=1$, although the idea can be applied to the general case.

4. The Complex Setting and System Form. In the remainder of this paper, we assume

1. $F : \mathbf{D} \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ can be extended to an analytic function in \mathbb{C}^n .

2. $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n) = ([\underline{x}_1, \bar{x}_1], \dots, [\underline{x}_n, \bar{x}_n])$ is a small box that will be constructed centered at the approximate solution \check{x} , i.e. $m(\mathbf{x}) = (\check{x}_1, \dots, \check{x}_n)$.
3. \check{x} is near a point x^* with $F(x^*) = 0$, such that $\|\check{x} - x^*\|$ is much smaller than the width of the box \mathbf{x} , and width of the box \mathbf{x} is small enough for a quadratic model to be accurate over the box \mathbf{x} .
4. F has been preconditioned as in Figure 3.1, and $F'(x^*)$ has null space of dimension 1.

So,

$$f_k(x) = (x_k - \check{x}_k) + \frac{\partial f_k}{\partial x_n}(\check{x})(x_n - \check{x}_n) + \mathcal{O}(\|x - \check{x}\|^2)$$

$$\text{for } 1 \leq k \leq n-1.$$

$$f_n(x) = \frac{1}{2} \sum_{k=1}^n \sum_{l=1}^n \frac{\partial^2 f_n}{\partial x_k \partial x_l}(\check{x})(x_k - \check{x}_k)(x_l - \check{x}_l) + \mathcal{O}(\|x - \check{x}\|^3)$$

For $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$, extend F to complex space: $x + iy$, with y in a small box

$$\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_n) = ([\underline{y}_1, \bar{y}_1], \dots, [\underline{y}_n, \bar{y}_n]),$$

where \mathbf{y} is centered at $(0, \dots, 0)$. Define $\mathbf{z} \equiv (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_n, \mathbf{y}_n) = ([\underline{x}_1, \bar{x}_1], [\underline{y}_1, \bar{y}_1], \dots, [\underline{x}_n, \bar{x}_n], [\underline{y}_n, \bar{y}_n])$, $u_k(x, y) \equiv \Re(f_k(x + iy))$ and $v_k(x, y) \equiv \Im(f_k(x + iy))$. With this, define

$$\tilde{F}(x, y) \equiv (u_1(x, y), v_1(x, y), \dots, u_n(x, y), v_n(x, y)) : \mathbb{R}^{2n} \rightarrow \mathbb{R}^{2n}.$$

Then, if preconditioning based on complete factorization of the midpoint matrix for $\mathbf{F}'(\mathbf{x})$ is used, the first-order terms are eliminated in the pattern of Figure 3.1, and, for $1 \leq k \leq (n-1)$,

$$(4.1) \quad \left. \begin{aligned} u_k(x, y) &= (x_k - \check{x}_k) + \frac{\partial f_k}{\partial x_n}(\check{x})(x_n - \check{x}_n) \\ &\quad + \mathcal{O}(\|(x - \check{x}, y)\|^2), \\ v_k(x, y) &= y_k + \frac{\partial f_k}{\partial x_n}(\check{x})y_n + \mathcal{O}(\|(x - \check{x}, y)\|^2), \end{aligned} \right\}$$

and

$$(4.2) \quad \left. \begin{aligned} u_n(x, y) &= \frac{1}{2} \sum_{k=1}^n \sum_{l=1}^n \frac{\partial^2 f_n}{\partial x_k \partial x_l}(\check{x})(x_k - \check{x}_k)(x_l - \check{x}_l) \\ &\quad - \frac{1}{2} \sum_{k=1}^n \sum_{l=1}^n \frac{\partial^2 f_n}{\partial x_k \partial x_l}(\check{x})y_k y_l \\ &\quad + \mathcal{O}(\|(x - \check{x}, y)\|^3), \\ v_n(x, y) &= \sum_{k=1}^n \sum_{l=1}^n \frac{\partial^2 f_n}{\partial x_k \partial x_l}(\check{x})(x_k - \check{x}_k)y_l \\ &\quad + \mathcal{O}(\|(x - \check{x}, y)\|^3). \end{aligned} \right\}$$

5. Simplification of a Degree Computation Procedure. To use Theorem 2.4 to compute the topological degree $d(\tilde{F}, \mathbf{z}, 0)$ directly in a verification algorithm would require a global search of the $4n(n-1)$ -dimensional faces of the $2n$ -dimensional box \mathbf{z} for zeros of \tilde{F}_{-p} . This is an inordinate amount of work for a verification process that would normally require only a single step of an interval Newton method in the nonsingular case. However, if the system is preconditioned and in the form described in §3 and §4, the verification can be reduced to $4n-4$ interval evaluations and four 1-dimensional searches.

To describe the simplification, define

$$\mathbf{x}_{\underline{k}} \equiv (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_{k-1}, \mathbf{y}_{k-1}, \underline{x}_k, \mathbf{y}_k, \mathbf{x}_{k+1}, \mathbf{y}_{k+1}, \dots, \mathbf{x}_n, \mathbf{y}_n),$$

and

$$\mathbf{x}_{\bar{k}} \equiv (\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_{k-1}, \mathbf{y}_{k-1}, \bar{x}_k, \mathbf{y}_k, \mathbf{x}_{k+1}, \mathbf{y}_{k+1}, \dots, \mathbf{x}_n, \mathbf{y}_n).$$

Similarly define $\mathbf{y}_{\underline{k}}$ and $\mathbf{y}_{\bar{k}}$. Also define

$$\tilde{F}_{-u_n}(x, y) \equiv (u_1(x, y), v_1(x, y), \dots, u_{n-1}(x, y), v_{n-1}(x, y), v_n(x, y)).$$

To compute the degree $d(\tilde{F}, \mathbf{z}, 0)$, we will consider \tilde{F}_{-u_n} on the boundary of \mathbf{z} . The boundary of \mathbf{z} consists of the $4n$ faces $\mathbf{x}_{\underline{1}}, \mathbf{x}_{\bar{1}}, \mathbf{y}_{\underline{1}}, \mathbf{y}_{\bar{1}}, \dots, \mathbf{x}_{\underline{n}}, \mathbf{x}_{\bar{n}}, \mathbf{y}_{\underline{n}}, \mathbf{y}_{\bar{n}}$.

Observe that, for $1 \leq k \leq (n-1)$,

$$\begin{aligned} \tilde{F}_{-u_n}(x, y) = 0 \text{ on } \mathbf{x}_{\underline{k}} \\ \implies u_k(x, y) &\approx (\underline{x}_k - \check{x}_k) + \frac{\partial f_k}{\partial x_n}(\check{x})(x_n - \check{x}_n) = 0 \\ \implies |\underline{x}_k - \check{x}_k| &= |\partial f_k / \partial x_n(\check{x})| |x_n - \check{x}_n| \\ \implies w(\mathbf{x}_k) &\leq |\partial f_k / \partial x_n(\check{x})| w(\mathbf{x}_n), \text{ or } \frac{w(\mathbf{x}_k)}{|\partial f_k / \partial x_n(\check{x})|} \leq w(\mathbf{x}_n). \end{aligned}$$

Similarly,

$$\begin{aligned} \tilde{F}_{-u_n}(x, y) = 0 \text{ on } \mathbf{x}_{\bar{k}} \\ \implies w(\mathbf{x}_k) &\leq |\partial f_k / \partial x_n(\check{x})| w(\mathbf{x}_n), \text{ or } \frac{w(\mathbf{x}_k)}{|\partial f_k / \partial x_n(\check{x})|} \leq w(\mathbf{x}_n). \end{aligned}$$

Also observe that, for $1 \leq k \leq (n-1)$,

$$\begin{aligned} \tilde{F}_{-u_n}(x, y) = 0 \text{ on } \mathbf{y}_{\underline{k}} \\ \implies u_k(x, y) &\approx \underline{y}_k + \frac{\partial f_k}{\partial x_n}(\check{x}) y_n = 0 \\ \implies |\underline{y}_k| &= |\partial f_k / \partial x_n(\check{x})| |y_n| \\ \implies w(\mathbf{y}_k) &\leq |\partial f_k / \partial x_n(\check{x})| w(\mathbf{y}_n), \text{ or } \frac{w(\mathbf{y}_k)}{|\partial f_k / \partial x_n(\check{x})|} \leq w(\mathbf{y}_n). \end{aligned}$$

Similarly,

$$\begin{aligned} \tilde{F}_{-u_n}(x, y) = 0 \text{ on } \mathbf{y}_{\bar{k}} \\ \implies w(\mathbf{y}_k) &\leq |\partial f_k / \partial x_n(\check{x})| w(\mathbf{y}_n), \text{ or } \frac{w(\mathbf{y}_k)}{|\partial f_k / \partial x_n(\check{x})|} \leq w(\mathbf{y}_n). \end{aligned}$$

Thus, if \mathbf{x}_n is chosen so that

$$(5.1) \quad w(\mathbf{x}_n) \leq \frac{1}{2} \min_{1 \leq k \leq n-1} \left\{ \frac{w(\mathbf{x}_k)}{|\partial f_k / \partial x_n(\check{x})|} \right\},$$

then it is unlikely that $u_k(x, y) = 0$ on either $\mathbf{x}_{\underline{k}}$ or $\mathbf{x}_{\bar{k}}$.

Similarly, if \mathbf{y}_n is chosen so that

$$(5.2) \quad w(\mathbf{y}_n) \leq \frac{1}{2} \min_{1 \leq k \leq n-1} \left\{ \frac{w(\mathbf{y}_k)}{|\partial f_k / \partial x_n(\check{x})|} \right\},$$

then it is unlikely that $v_k(x, y) = 0$ on either $\mathbf{y}_{\underline{k}}$ or $\mathbf{y}_{\bar{k}}$.

Here, the coefficient $\frac{1}{2}$ is to take into consideration the fact that $u_k(x, y) \approx (x_k - \tilde{x}_k) + \frac{\partial f_k}{\partial x_n}(\tilde{x})(x_n - \tilde{x}_n)$ and $v_k(x, y) \approx y_k + \frac{\partial f_k}{\partial x_n}(\tilde{x})y_n$ are only approximate equalities.

When $\partial f_k/\partial x_n(\tilde{x})$ happens to be 0, we can think that $\frac{w(\mathbf{x}_k)}{|\partial f_k/\partial x_n(\tilde{x})|}$ and $\frac{w(\mathbf{y}_k)}{|\partial f_k/\partial x_n(\tilde{x})|}$ are $+\infty$.

By constructing the box \mathbf{z} in this way, we can eliminate search of $4n - 4$ of the $4n$ faces of the boundary of \mathbf{z} , since we have arranged to verify $\tilde{F}_{-u_n}(x, y) \neq 0$ on each of these faces. Elimination of these $4n - 4$ faces only needs $4n - 4$ interval evaluations. Then, we only need to search the four faces $\mathbf{x}_{\underline{n}}, \mathbf{x}_{\bar{n}}, \mathbf{y}_{\underline{n}}$ and $\mathbf{y}_{\bar{n}}$ for solutions of $\tilde{F}_{-u_n}(x, y) = 0$, regardless of how large n is. This reduces total computational cost dramatically, since searching a face is expensive. The following theorem forms the theoretical basis of our algorithm in §6.1.

THEOREM 5.1. *Suppose*

1. $u_k \neq 0$ on $\mathbf{x}_{\underline{k}}$ and $\mathbf{x}_{\bar{k}}$, and $v_k \neq 0$ on $\mathbf{y}_{\underline{k}}$ and $\mathbf{y}_{\bar{k}}$, $k = 1, \dots, n - 1$;
2. $\tilde{F}_{-u_n} = 0$ has a unique solution on $\mathbf{x}_{\underline{n}}$ and $\mathbf{x}_{\bar{n}}$ with y_n in the interior of $\mathbf{y}_{\underline{n}}$, and $\tilde{F}_{-u_n} = 0$ has a unique solution on $\mathbf{y}_{\underline{n}}$ and $\mathbf{y}_{\bar{n}}$ with x_n in the interior of $\mathbf{x}_{\underline{n}}$;
3. $u_n \neq 0$ at the four solutions of $\tilde{F}_{-u_n} = 0$ in condition 2; and
4. the Jacobi matrices of \tilde{F}_{-u_n} are non-singular at the four solutions of $\tilde{F}_{-u_n} = 0$ in condition 2.

Then

$$\begin{aligned} d(\tilde{F}, \mathbf{z}, 0) = & - \sum_{\substack{x_n = \underline{x}_n \\ \tilde{F}_{-u_n}(x, y) = 0 \\ u_n(x, y) > 0}} \operatorname{sgn} \left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} y_n}(x, y) \right| \\ & + \sum_{\substack{x_n = \bar{x}_n \\ \tilde{F}_{-u_n}(x, y) = 0 \\ u_n(x, y) > 0}} \operatorname{sgn} \left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} y_n}(x, y) \right| \\ & + \sum_{\substack{y_n = \underline{y}_n \\ \tilde{F}_{-u_n}(x, y) = 0 \\ u_n(x, y) > 0}} \operatorname{sgn} \left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} x_n}(x, y) \right| \\ & - \sum_{\substack{y_n = \bar{y}_n \\ \tilde{F}_{-u_n}(x, y) = 0 \\ u_n(x, y) > 0}} \operatorname{sgn} \left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} x_n}(x, y) \right|. \end{aligned}$$

Proof. Condition 1 implies $\tilde{F} \neq 0$ on $\mathbf{x}_{\underline{k}}, \mathbf{x}_{\bar{k}}, \mathbf{y}_{\underline{k}}$ and $\mathbf{y}_{\bar{k}}$, $k = 1, \dots, n - 1$, and 2 and 3 imply $\tilde{F} \neq 0$ on $\mathbf{x}_{\underline{n}}, \mathbf{x}_{\bar{n}}, \mathbf{y}_{\underline{n}}$ and $\mathbf{y}_{\bar{n}}$. So, $\tilde{F} \neq 0$ on $\partial \mathbf{z}$.

Condition 1 implies $\tilde{F}_{-u_n} \neq 0$ on $\partial \mathbf{x}_{\underline{k}}, \partial \mathbf{x}_{\bar{k}}, \partial \mathbf{y}_{\underline{k}}$ and $\partial \mathbf{y}_{\bar{k}}$, $k = 1, \dots, n - 1$. $\partial \mathbf{x}_{\underline{n}}$ consists of $2(n - 1)$ $(n - 2)$ -dimensional boxes each of which is either embedded in some $\mathbf{x}_{\underline{k}}, \mathbf{x}_{\bar{k}}, \mathbf{y}_{\underline{k}}$ or $\mathbf{y}_{\bar{k}}$, $1 \leq k \leq n - 1$, or is embedded in $\mathbf{y}_{\underline{n}}$ or $\mathbf{y}_{\bar{n}}$. So, by 2 and 3, $\tilde{F}_{-u_n} \neq 0$ on $\partial \mathbf{x}_{\underline{n}}$. Similarly, $\tilde{F}_{-u_n} \neq 0$ on $\partial \mathbf{x}_{\bar{n}}, \partial \mathbf{y}_{\underline{n}}$ and $\partial \mathbf{y}_{\bar{n}}$. So, condition 1 in Theorem 2.4 is satisfied.

Condition 2 in Theorem 2.4 is automatically satisfied since $F_{-p} = 0$ either has no solutions or a unique solution on $\mathbf{x}_k, \mathbf{x}_{\bar{k}}, \mathbf{y}_k$, or $\mathbf{y}_{\bar{k}}, 1 \leq k \leq n$.

Then, by noticing 4, the conditions of Theorem 2.4 are satisfied. The formula is thus obtained with $s = +1$. \square

The conditions of Theorem 5.1 will be satisfied when the system is that as described in §3 and §4, the box \mathbf{z} is constructed as in (5.1) and (5.2), and the quadratic model is accurate. (See Theorem 5.2 and its proof of the results when all the approximations are exact.)

In Theorem 5.1, the degree consists of contributions of the four faces we search. We can compute the degree contribution of each of the four faces and then add them to get the degree.

In Theorem 5.1 we choose $s = +1$. We can also choose $s = -1$. That doesn't make any difference in our context if we ignore higher order terms in the values of u_n at the solutions of $\tilde{F}_{-u_n} = 0$ on the four faces $\mathbf{x}_n, \mathbf{x}_{\bar{n}}, \mathbf{y}_n$ and $\mathbf{y}_{\bar{n}}$. To be specific, the four values of u_n are

$$\begin{aligned} u_n &= \frac{1}{2}\Delta(\underline{x}_n - \check{x}_n)^2 + \mathcal{O}\left(\|(x - \check{x}, y)\|^3\right), \\ u_n &= \frac{1}{2}\Delta(\bar{x}_n - \check{x}_n)^2 + \mathcal{O}\left(\|(x - \check{x}, y)\|^3\right), \\ u_n &= -\frac{1}{2}\Delta\underline{y}_n^2 + \mathcal{O}\left(\|(x - \check{x}, y)\|^3\right), \\ u_n &= -\frac{1}{2}\Delta\bar{y}_n^2 + \mathcal{O}\left(\|(x - \check{x}, y)\|^3\right), \end{aligned}$$

respectively, where, Δ is defined in (5.3). When we choose $w(\mathbf{y}_k)$ the same (or roughly the same) as $w(\mathbf{x}_k)$, the values of u_n as a function of \underline{y}_n (or \bar{y}_n) will be the same (or roughly the same) as the values of u_n as a function of $\underline{x}_n - \check{x}_n$ (or $\bar{x}_n - \check{x}_n$). So, if we ignore higher order terms, the cost of verifying $u_n < 0$ and searching for solutions of $\tilde{F}_{-u_n} = 0$ with $u_n > 0$ is approximately the same as the cost of verifying $u_n > 0$ and searching for solutions of $\tilde{F}_{-u_n} = 0$ with $u_n < 0$.

Next, we will give a theorem which will further reduce the search cost by telling us how we should search. Before introducing the theorem, define

$$\begin{aligned} \alpha_k &\equiv \frac{\partial f_k}{\partial x_n}(\check{x}), & 1 \leq k \leq n-1, \\ \alpha_n &\equiv -1, \\ \beta_{kl} &\equiv \frac{\partial^2 f_n}{\partial x_k \partial x_l}(\check{x}) & 1 \leq k \leq n, 1 \leq l \leq n, \\ (5.3) \quad \Delta &\equiv \sum_{k=1}^n \sum_{l=1}^n \beta_{kl} \alpha_k \alpha_l \end{aligned}$$

Note $\beta_{kl} = \beta_{lk}$

THEOREM 5.2. *If the approximations of (4.1) and (4.2) are exact, if we construct the box \mathbf{z} as in (5.1) and (5.2), and if $\Delta \neq 0$, then $d(\tilde{F}, \mathbf{z}, 0) = 2$.*

Proof. Under the assumptions,

$$(5.4) \quad u_k = (x_k - \check{x}_k) + \alpha_k(x_n - \check{x}_n), \quad 1 \leq k \leq n-1,$$

$$(5.5) \quad v_k = y_k + \alpha_k y_n, \quad 1 \leq k \leq n-1,$$

$$(5.6) \quad u_n = \frac{1}{2} \sum_{k=1}^n \sum_{l=1}^n \beta_{kl} (x_k - \check{x}_k)(x_l - \check{x}_l) - \frac{1}{2} \sum_{k=1}^n \sum_{l=1}^n \beta_{kl} y_k y_l$$

$$(5.7) \quad v_n = \sum_{k=1}^n \sum_{l=1}^n \beta_{kl} (x_k - \check{x}_k) y_l$$

Due to the construction of the box \mathbf{z} ,

$$u_k = (\underline{x}_k - \check{x}_k) + \alpha_k (x_n - \check{x}_n) \neq 0$$

on $\underline{\mathbf{x}}_k$ and $\underline{\mathbf{x}}_{\bar{k}}$, and

$$v_k = \underline{y}_k + \alpha_k y_n \neq 0$$

on $\underline{\mathbf{y}}_k$ and $\underline{\mathbf{y}}_{\bar{k}}$, where $k = 1, \dots, n-1$.

Next, we locate the solutions of $\tilde{F}_{-u_n} = 0$ on $\underline{\mathbf{x}}_n$, $\underline{\mathbf{x}}_{\bar{n}}$, $\underline{\mathbf{y}}_n$ and $\underline{\mathbf{y}}_{\bar{n}}$.

1) On $\underline{\mathbf{x}}_n$;

$$(5.8) \quad u_k = 0 \implies \tilde{x}_k = \check{x}_k - \alpha_k (\underline{x}_n - \check{x}_n), \quad 1 \leq k \leq n-1,$$

$$(5.9) \quad v_k = 0 \implies \tilde{y}_k = -\alpha_k y_n, \quad 1 \leq k \leq n-1,$$

Plugging equations (5.8) and (5.9) into equations (5.6) and (5.7), we get

$$(5.10) \quad u_n = \frac{1}{2} \sum_{k=1}^n \sum_{l=1}^n \beta_{kl} \alpha_k \alpha_l (\underline{x}_n - \check{x}_n)^2 - \frac{1}{2} \sum_{k=1}^n \sum_{l=1}^n \beta_{kl} \alpha_k \alpha_l y_n^2$$

$$= \frac{1}{2} \Delta (\underline{x}_n - \check{x}_n)^2 - \frac{1}{2} \Delta y_n^2$$

$$(5.11) \quad v_n = \sum_{k=1}^n \sum_{l=1}^n \beta_{kl} \alpha_k \alpha_l (\underline{x}_n - \check{x}_n) y_n$$

$$= \Delta (\underline{x}_n - \check{x}_n) y_n.$$

Then,

$$(5.12) \quad v_n = 0 \implies \tilde{y}_n = 0,$$

since $\Delta \neq 0$.

Thus, by equation (5.9)

$$(5.13) \quad \tilde{y}_n = 0 \implies \tilde{y}_k = 0, \quad 1 \leq k \leq n-1,$$

So, $\tilde{F}_{-u_n} = 0$ has a unique solution

$$(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = (\tilde{x}_1, 0, \dots, \tilde{x}_{n-1}, 0, \underline{x}_n, 0)$$

on $\underline{\mathbf{x}}_n$.

Plugging (5.12) into (5.10), we get the u_n value at this solution, which is

$$(5.14) \quad u_n = \frac{1}{2} \Delta(\underline{x}_n - \check{x}_n)^2.$$

Next, we compute the determinant of the Jacobi matrix of \tilde{F}_{-u_n} at this solution. Define

$$\gamma_k \equiv \sum_{l=1}^n \beta_{lk} \alpha_l.$$

Noting equations (5.4), (5.5) and (5.7), we have

$$(5.15) \quad \begin{aligned} & \left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} y_n}(\tilde{x}, \tilde{y}) \right| \\ &= \begin{vmatrix} \frac{\partial u_1}{\partial x_1} & \frac{\partial u_1}{\partial y_1} & \cdots & \frac{\partial u_1}{\partial x_{n-1}} & \frac{\partial u_1}{\partial y_{n-1}} & \frac{\partial u_1}{\partial y_n} \\ \frac{\partial v_1}{\partial x_1} & \frac{\partial v_1}{\partial y_1} & \cdots & \frac{\partial v_1}{\partial x_{n-1}} & \frac{\partial v_1}{\partial y_{n-1}} & \frac{\partial v_1}{\partial y_n} \\ \vdots & \vdots & & \vdots & \vdots & \vdots \\ \frac{\partial u_{n-1}}{\partial x_1} & \frac{\partial u_{n-1}}{\partial y_1} & \cdots & \frac{\partial u_{n-1}}{\partial x_{n-1}} & \frac{\partial u_{n-1}}{\partial y_{n-1}} & \frac{\partial u_{n-1}}{\partial y_n} \\ \frac{\partial v_{n-1}}{\partial x_1} & \frac{\partial v_{n-1}}{\partial y_1} & \cdots & \frac{\partial v_{n-1}}{\partial x_{n-1}} & \frac{\partial v_{n-1}}{\partial y_{n-1}} & \frac{\partial v_{n-1}}{\partial y_n} \\ \frac{\partial v_n}{\partial x_1} & \frac{\partial v_n}{\partial y_1} & \cdots & \frac{\partial v_n}{\partial x_{n-1}} & \frac{\partial v_n}{\partial y_{n-1}} & \frac{\partial v_n}{\partial y_n} \end{vmatrix} \\ &= \begin{vmatrix} 1 & 0 & \cdots & 0 & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 & \alpha_1 \\ \vdots & \vdots & & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 & 0 \\ 0 & 0 & \cdots & 0 & 1 & \alpha_{n-1} \\ 0 & -\gamma_1(\underline{x}_n - \check{x}_n) & \cdots & 0 & -\gamma_{n-1}(\underline{x}_n - \check{x}_n) & -\gamma_n(\underline{x}_n - \check{x}_n) \end{vmatrix} \\ &= -(\underline{x}_n - \check{x}_n) \begin{vmatrix} 1 & 0 & \cdots & 0 & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 & \alpha_1 \\ \vdots & \vdots & & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 & 0 \\ 0 & 0 & \cdots & 0 & 1 & \alpha_{n-1} \\ 0 & \gamma_1 & \cdots & 0 & \gamma_{n-1} & \gamma_n \end{vmatrix} \\ &= -(\underline{x}_n - \check{x}_n) \left(-\sum_{k=1}^n \alpha_k \gamma_k \right) \\ &= (\underline{x}_n - \check{x}_n) \sum_{k=1}^n \alpha_k \sum_{l=1}^n \beta_{lk} \alpha_l \\ &= (\underline{x}_n - \check{x}_n) \sum_{k=1}^n \sum_{l=1}^n \beta_{lk} \alpha_k \alpha_l \\ &= (\underline{x}_n - \check{x}_n) \sum_{k=1}^n \sum_{l=1}^n \beta_{kl} \alpha_k \alpha_l \\ &= (\underline{x}_n - \check{x}_n) \Delta. \end{aligned}$$

2) On $\mathbf{x}_{\bar{n}}$;

Similarly, $\tilde{F}_{\neg u_n} = 0$ has a unique solution (\tilde{x}, \tilde{y}) on $\mathbf{x}_{\bar{n}}$.

The u_n value at this solution is

$$(5.16) \quad u_n = \frac{1}{2} \Delta (\bar{x}_n - \check{x}_n)^2.$$

The determinant of the Jacobi matrix of $\tilde{F}_{\neg u_n}$ at this solution is

$$(5.17) \quad \left| \frac{\partial \tilde{F}_{\neg u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} y_n}(\tilde{x}, \tilde{y}) \right| = (\bar{x}_n - \check{x}_n) \Delta$$

3) On $\mathbf{y}_{\underline{n}}$;

$$(5.18) \quad u_k = 0 \implies \tilde{x}_k = \check{x}_k - \alpha_k (x_n - \check{x}_n), \quad 1 \leq k \leq n-1,$$

$$(5.19) \quad v_k = 0 \implies \tilde{y}_k = -\alpha_k \underline{y}_n, \quad 1 \leq k \leq n-1,$$

Plugging equations (5.18) and (5.19) into equations (5.6) and (5.7), we get

$$(5.20) \quad \begin{aligned} u_n &= \frac{1}{2} \sum_{k=1}^n \sum_{l=1}^n \beta_{kl} \alpha_k \alpha_l (x_n - \check{x}_n)^2 - \frac{1}{2} \sum_{k=1}^n \sum_{l=1}^n \beta_{kl} \alpha_k \alpha_l \underline{y}_n^2 \\ &= \frac{1}{2} \Delta (x_n - \check{x}_n)^2 - \frac{1}{2} \Delta \underline{y}_n^2 \end{aligned}$$

$$(5.21) \quad \begin{aligned} v_n &= \sum_{k=1}^n \sum_{l=1}^n \beta_{kl} \alpha_k \alpha_l (x_n - \check{x}_n) \underline{y}_n \\ &= \Delta (x_n - \check{x}_n) \underline{y}_n. \end{aligned}$$

Then,

$$(5.22) \quad v_n = 0 \implies \tilde{x}_n = \check{x}_n,$$

since $\Delta \neq 0$.

Thus, by equation (5.18)

$$(5.23) \quad \tilde{x}_n = \check{x}_n \implies \tilde{x}_k = \check{x}_k, \quad 1 \leq k \leq n-1,$$

So, $\tilde{F}_{\neg u_n} = 0$ has a unique solution

$$(\tilde{x}, \tilde{y}) = (\check{x}_1, \check{y}_1, \dots, \check{x}_{n-1}, \check{y}_{n-1}, \check{x}_n, \underline{y}_n)$$

on $\mathbf{y}_{\underline{n}}$.

Plugging (5.22) into (5.20), we get the u_n value at this solution, which is

$$(5.24) \quad u_n = -\frac{1}{2} \Delta \underline{y}_n^2.$$

Next, we compute the determinant of the Jacobi matrix of $\tilde{F}_{\neg u_n}$ at this solution.

Noting the equations (5.4), (5.5) and (5.7), we have

$$\begin{aligned}
(5.25) \quad & \left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} x_n}(\tilde{x}, \tilde{y}) \right| \\
= & \begin{vmatrix} \frac{\partial u_1}{\partial x_1} & \frac{\partial u_1}{\partial y_1} & \cdots & \frac{\partial u_1}{\partial x_{n-1}} & \frac{\partial u_1}{\partial y_{n-1}} & \frac{\partial u_1}{\partial x_n} \\ \frac{\partial v_1}{\partial x_1} & \frac{\partial v_1}{\partial y_1} & \cdots & \frac{\partial v_1}{\partial x_{n-1}} & \frac{\partial v_1}{\partial y_{n-1}} & \frac{\partial v_1}{\partial x_n} \\ \vdots & \vdots & & \vdots & \vdots & \vdots \\ \frac{\partial u_{n-1}}{\partial x_1} & \frac{\partial u_{n-1}}{\partial y_1} & \cdots & \frac{\partial u_{n-1}}{\partial x_{n-1}} & \frac{\partial u_{n-1}}{\partial y_{n-1}} & \frac{\partial u_{n-1}}{\partial x_n} \\ \frac{\partial v_{n-1}}{\partial x_1} & \frac{\partial v_{n-1}}{\partial y_1} & \cdots & \frac{\partial v_{n-1}}{\partial x_{n-1}} & \frac{\partial v_{n-1}}{\partial y_{n-1}} & \frac{\partial v_{n-1}}{\partial x_n} \\ \frac{\partial v_n}{\partial x_1} & \frac{\partial v_n}{\partial y_1} & \cdots & \frac{\partial v_n}{\partial x_{n-1}} & \frac{\partial v_n}{\partial y_{n-1}} & \frac{\partial v_n}{\partial x_n} \end{vmatrix} \\
= & \begin{vmatrix} 1 & 0 & \cdots & 0 & 0 & \alpha_1 \\ 0 & 1 & \cdots & 0 & 0 & 0 \\ \vdots & \vdots & & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 & \alpha_{n-1} \\ 0 & 0 & \cdots & 0 & 1 & 0 \\ -\gamma_1 \underline{y}_n & 0 & \cdots & -\gamma_{n-1} \underline{y}_n & 0 & -\gamma_n \underline{y}_n \end{vmatrix} \\
= & -\underline{y}_n \begin{vmatrix} 1 & 0 & \cdots & 0 & 0 & \alpha_1 \\ 0 & 1 & \cdots & 0 & 0 & 0 \\ \vdots & \vdots & & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 & \alpha_{n-1} \\ 0 & 0 & \cdots & 0 & 1 & 0 \\ \gamma_1 & 0 & \cdots & \gamma_{n-1} & 0 & \gamma_n \end{vmatrix} \\
= & -\underline{y}_n \left(-\sum_{k=1}^n \alpha_k \gamma_k \right) \\
= & \underline{y}_n \sum_{k=1}^n \alpha_k \sum_{l=1}^n \beta_{lk} \alpha_l \\
= & \underline{y}_n \sum_{k=1}^n \sum_{l=1}^n \beta_{lk} \alpha_k \alpha_l \\
= & \underline{y}_n \sum_{k=1}^n \sum_{l=1}^n \beta_{kl} \alpha_k \alpha_l \\
= & \underline{y}_n \Delta.
\end{aligned}$$

4) On $\mathbf{y}_{\bar{n}}$;

Similarly, $\tilde{F}_{-u_n} = 0$ has a unique solution (\tilde{x}, \tilde{y}) on $\mathbf{y}_{\bar{n}}$.

The u_n value at this solution is

$$(5.26) \quad u_n = -\frac{1}{2} \Delta \bar{y}_n^2.$$

The determinant of the Jacobi matrix of \tilde{F}_{-u_n} at this solution is

$$(5.27) \quad \left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} x_n}(\tilde{x}, \tilde{y}) \right| = \bar{y}_n \Delta$$

Finally, we can use the formula in Theorem 5.1 to compute the topological degree $d(\tilde{F}, \mathbf{z}, 0)$.

If $\Delta > 0$, then we know from the equations (5.14), (5.16), (5.24) and (5.26) that u_n is only positive at the solutions of $\tilde{F}_{-u_n} = 0$ on \mathbf{x}_n and $\mathbf{x}_{\bar{n}}$. We also know the signs of the determinants of the Jacobi matrices at the two solutions from equations (5.15) and (5.17). So,

$$d(\tilde{F}, \mathbf{z}, 0) = -(-1) + (+1) = 2$$

If $\Delta < 0$, then we know from the equations (5.14), (5.16), (5.24) and (5.26) that u_n is only positive at the solutions of $\tilde{F}_{-u_n} = 0$ on \mathbf{y}_n and $\mathbf{y}_{\bar{n}}$. We also know the signs of the determinants of the Jacobi matrices at the two solutions from equations (5.25) and (5.27). So,

$$d(\tilde{F}, \mathbf{z}, 0) = +(+1) - (-1) = 2$$

□

The proof of Theorem 5.2 tells us approximately where we can expect to find the solutions of $\tilde{F}_{-u_n} = 0$ on the four faces we search and the value of the degree we can expect when the approximations (4.1) and (4.2) are accurate.

From (4.1), we know that if x_n is known precisely, formally solving $\mathbf{u}_k(\mathbf{x}, \mathbf{y}) = 0$ for x_k gives sharper bounds $\tilde{\mathbf{x}}_k$ with $w(\tilde{\mathbf{x}}_k) = \mathcal{O}(\|(\mathbf{x} - \tilde{\mathbf{x}}, \mathbf{y})\|^2)$, $1 \leq k \leq n - 1$. Similarly, if y_n is known precisely, formally solving $\mathbf{v}_k(\mathbf{x}, \mathbf{y}) = 0$ for y_k gives sharper bounds $\tilde{\mathbf{y}}_k$ with $w(\tilde{\mathbf{y}}_k) = \mathcal{O}(\|(\mathbf{x} - \tilde{\mathbf{x}}, \mathbf{y})\|^2)$, $1 \leq k \leq n - 1$. So, when we search \mathbf{x}_n (or $\mathbf{x}_{\bar{n}}$) for solutions of $\tilde{F}_{-u_n} = 0$, we can first get sharper bounds for x_k , $1 \leq k \leq n - 1$, since x_n is known precisely. Then, for a small subinterval \mathbf{y}_n^0 of \mathbf{y}_n , we can solve $\mathbf{v}_k(\mathbf{x}, \mathbf{y}) = 0$ for y_k to get sharper bounds $\tilde{\mathbf{y}}_k$ with $w(\tilde{\mathbf{y}}_k) = \mathcal{O}(\max(\|(\mathbf{x} - \tilde{\mathbf{x}}, \mathbf{y})\|^2, \|\mathbf{y}_n^0\|))$, $1 \leq k \leq n - 1$. Thus we get a small subface of \mathbf{x}_n (or $\mathbf{x}_{\bar{n}}$) over which we can either use an interval Newton method to verify the existence and uniqueness of the zero of \tilde{F}_{-u_n} or use mean-value extensions to verify that \tilde{F}_{-u_n} has no zeros, depending on whether \mathbf{y}_n^0 is in the middle of \mathbf{y}_n or not. So, we end up with searching over a 1-dimensional interval \mathbf{y}_n . This further reduces the search cost. We can search \mathbf{y}_n or $\mathbf{y}_{\bar{n}}$ in a similar way.

The analysis in this section leads to a practical algorithm in the next section.

6. The Algorithm and Its Computational Complexity. In this section, we will present the actual algorithm and give its complexity.

6.1. Algorithm. The algorithm consists of three phases. In the box-setting phase, we set the box \mathbf{z} . In the elimination phase, we verify that $u_k \neq 0$ on \mathbf{x}_k and $\mathbf{x}_{\bar{k}}$, and $v_k \neq 0$ on \mathbf{y}_k and $\mathbf{y}_{\bar{k}}$, where $1 \leq k \leq n - 1$. In the search phase, we verify the unique solution of $\tilde{F}_{-u_n} = 0$ on \mathbf{x}_n and $\mathbf{x}_{\bar{n}}$ with y_n in the interior of \mathbf{y}_n , and on \mathbf{y}_n and $\mathbf{y}_{\bar{n}}$ with x_n in the interior of \mathbf{x}_n , compute the signs of u_n and the signs of the Jacobi matrices of \tilde{F}_{-u_n} at the four solutions of $\tilde{F}_{-u_n} = 0$, compute the degree

contributions of the 4 faces \mathbf{x}_n , $\mathbf{x}_{\bar{n}}$, \mathbf{y}_n and $\mathbf{y}_{\bar{n}}$ according to the formula in Theorem 5.1 and finally add the contributions to get the degree.

ALGORITHM

Box-setting Phase

1. Compute the preconditioner of the original system, using Gaussian elimination with full pivoting.
2. Set the widths of \mathbf{x}_k and \mathbf{y}_k (see explanation below), for $1 \leq k \leq n-1$.
3. Set the widths of \mathbf{x}_n and \mathbf{y}_n as in (5.1) and (5.2).

Elimination Phase

1. Do for $1 \leq k \leq n-1$
 - (a) Do for \mathbf{x}_k and $\mathbf{x}_{\bar{k}}$
 - i. Compute the mean-value extension of \mathbf{u}_k over that face.
 - ii. If $0 \in \mathbf{u}_k$, then, stop and signal failure.
 - (b) Do for \mathbf{y}_k and $\mathbf{y}_{\bar{k}}$
 - i. Compute the mean-value extension of \mathbf{v}_k over that face.
 - ii. If $0 \in \mathbf{v}_k$, then, stop and signal failure.

Search Phase

1. Do for \mathbf{x}_n and $\mathbf{x}_{\bar{n}}$
 - (a)
 - i. Use mean-value extensions for $\mathbf{u}_k(\mathbf{x}, \mathbf{y}) = 0$ to solve for x_k to get sharper bounds $\tilde{\mathbf{x}}_k$ with width $\mathcal{O}(\|(\mathbf{x} - \tilde{\mathbf{x}}, \mathbf{y})\|^2)$, $1 \leq k \leq n-1$.
 - ii. If $\tilde{\mathbf{x}}_k \cap \mathbf{x}_k = \emptyset$, then return the degree contribution of that face as 0.
 - iii. Update \mathbf{x}_k .
 - (b)
 - i. Compute the mean-value extension \mathbf{u}_n over that face.
 - ii. If $\mathbf{u}_n < 0$, then return the degree contribution of that face as 0.
 - (c) Construct a small subinterval \mathbf{y}_n^0 of \mathbf{y}_n which is centered at 0.
 - (d)
 - i. Use mean-value extensions for $\mathbf{v}_k(\mathbf{x}, \mathbf{y}) = 0$ to solve for y_k to get sharper bounds $\tilde{\mathbf{y}}_k$ with width $\mathcal{O}(\max(\|(\mathbf{x} - \tilde{\mathbf{x}}, \mathbf{y})\|^2, \|\mathbf{y}_n^0\|))$, $1 \leq k \leq n-1$, and thus to get a subface \mathbf{x}_n^0 (or $\mathbf{x}_{\bar{n}}^0$) of \mathbf{x}_n (or $\mathbf{x}_{\bar{n}}$).
 - ii. If $\tilde{\mathbf{y}}_k \cap \mathbf{y}_k = \emptyset$, then stop and signal failure.
 - (e)
 - i. Set up an interval Newton method for \tilde{F}_{-u_n} to verify existence and uniqueness of a zero in the subface \mathbf{x}_n^0 (or $\mathbf{x}_{\bar{n}}^0$).
 - ii. If the zero can not be verified, then stop and signal failure.
 - (f) Inflate \mathbf{y}_n^0 as much as possible provided the existence and uniqueness of the zero of \tilde{F}_{-u_n} can be verified over the corresponding subface, and thus get a subinterval \mathbf{y}_n^1 of \mathbf{y}_n .
 - (g) In this step, we verify $\tilde{F}_{-u_n} = 0$ has no solutions when $y_n \in \mathbf{y}_n \setminus \mathbf{y}_n^1$. $\mathbf{y}_n \setminus \mathbf{y}_n^1$ has two separate parts; we denote the lower part by \mathbf{y}_n^l and the upper part by \mathbf{y}_n^u . We only present the processing of the lower part. The upper part can be processed in a similar way.
 - i. Do
 - A. Use mean-value extensions for $\mathbf{v}_k(\mathbf{x}, \mathbf{y}) = 0$ to solve for y_k to get sharper bounds for y_k , $1 \leq k \leq n-1$, and thus to get a subface of \mathbf{x}_n (or $\mathbf{x}_{\bar{n}}$).
 - B. Compute the mean-value extensions \tilde{F}_{-u_n} over the subface obtained in the last step.
 - C. If $0 \in \tilde{F}_{-u_n}$, then bisection \mathbf{y}_n^l , update the lower part as a new \mathbf{y}_n^l .

and cycle.

If $0 \notin \tilde{F}_{-u_n}$, then exit the loop.

ii. Do

A. If $\underline{y}_n^1 \leq \bar{y}_n^l$, then exit the loop.

B. $\mathbf{y}_n^l \leftarrow [\bar{y}_n^l, \bar{y}_n^l + w(\mathbf{y}_n^l)]$.

C. Use mean-value extensions for $\mathbf{v}_k(\mathbf{x}, \mathbf{y}) = 0$ to solve for y_k to get sharper bounds for y_k , $1 \leq k \leq n-1$, and thus to get a subface of \mathbf{x}_n (or $\mathbf{x}_{\bar{n}}$.)

D. Compute the mean-value extensions \tilde{F}_{-u_n} over the subface obtained in the last step.

E. If $0 \notin \tilde{F}_{-u_n}$, then cycle.

If $0 \in \tilde{F}_{-u_n}$, then $\mathbf{y}_n^l \leftarrow [\underline{y}_n^l, \text{mid}(\mathbf{y}_n^l)]$ and cycle.

(h) i. Compute the mean-value extension of \mathbf{u}_n over \mathbf{x}_n^0 (or $\mathbf{x}_{\bar{n}}^0$.)

ii. If $\mathbf{u}_n < 0$, then return the degree contribution of that face as 0.

(i) i. Compute $\left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} y_n}(\mathbf{x}_n^0) \right|$ (or $\left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} y_n}(\mathbf{x}_{\bar{n}}^0) \right|$).

ii. If $0 \in \left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} y_n}(\mathbf{x}_n^0) \right|$ (or $0 \in \left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} y_n}(\mathbf{x}_{\bar{n}}^0) \right|$), then stop and signal failure.

(j) Use the formula in Theorem 5.1 to compute the degree contribution of that face.

2. Do for \mathbf{y}_n and $\mathbf{y}_{\bar{n}}$

(a) Same as Step 1a except change x_k to y_k , $\tilde{\mathbf{x}}_k$ to $\tilde{\mathbf{y}}_k$, \mathbf{x}_k to \mathbf{y}_k and \mathbf{u}_k to \mathbf{v}_k .

(b) Same as Step 1b.

(c) Same as Step 1c except change \mathbf{y}_n^0 to \mathbf{x}_n^0 , \mathbf{y}_n to \mathbf{x}_n and 0 to \tilde{x}_n .

(d) Same as Step 1d except change y_k to x_k , $\tilde{\mathbf{y}}_k$ to $\tilde{\mathbf{x}}_k$, \mathbf{y}_k to \mathbf{x}_k , \mathbf{x}_n^0 to \mathbf{y}_n^0 , $\mathbf{x}_{\bar{n}}^0$ to $\mathbf{y}_{\bar{n}}^0$, \mathbf{x}_n to \mathbf{y}_n and $\mathbf{x}_{\bar{n}}$ to $\mathbf{y}_{\bar{n}}$.

(e) Same as Step 1e except change \mathbf{x}_n^0 to \mathbf{y}_n^0 and $\mathbf{x}_{\bar{n}}^0$ to $\mathbf{y}_{\bar{n}}^0$.

(f) Same as Step 1f except change \mathbf{y}_n^0 to \mathbf{x}_n^0 , \mathbf{y}_n^1 to \mathbf{x}_n^1 and \mathbf{y}_n to \mathbf{x}_n .

(g) Same as Step 1g except change $\mathbf{y}_n \setminus \mathbf{y}_n^1$ to $\mathbf{x}_n \setminus \mathbf{x}_n^1$.

(h) Same as Step 1h except change \mathbf{x}_n^0 to \mathbf{y}_n^0 and $\mathbf{x}_{\bar{n}}^0$ to $\mathbf{y}_{\bar{n}}^0$.

(i) Same as Step 1i except change

$$\left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} y_n}(\mathbf{x}_n^0) \right| \text{ to } \left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} x_n}(\mathbf{y}_n^0) \right| \text{ and } \left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} y_n}(\mathbf{x}_{\bar{n}}^0) \right| \text{ to } \left| \frac{\partial \tilde{F}_{-u_n}}{\partial x_1 y_1 \dots x_{n-1} y_{n-1} x_n}(\mathbf{y}_{\bar{n}}^0) \right|.$$

(j) Same as Step 1j.

3. Add the degree contributions of the four faces obtained in steps 1 and 2 to get the degree.

END OF ALGORITHM

An Explanation of the Algorithm

1. In the box-setting phase, in Step 2, the width $w(\mathbf{x}_k)$ of \mathbf{x}_k depends on the accuracy of the approximate solution \tilde{x} of the system $F(\mathbf{x}) = 0$. $w(\mathbf{x}_k)$ should be much larger than $|\tilde{x}_k - x_k^*|$. But, at the same time, it should not be too large since the quadratic model needs to be accurate over the box.
2. In the search phase, in Step 1b (or 2b), we check the sign of u_n on that face and discard that face at the earliest possible time if $u_n < 0$ on that face, since we know the degree contribution of that face is 0 according to the formula

in Theorem 5.1. This will save time significantly if it happens that $u_n < 0$ on that face. It did happen for all the test problems. (See §8 for the test results.)

3. In the search phase, in Step 1e (or 2e), we precondition the system \tilde{F}_{-u_n} before we use an interval Newton method, so that the method will succeed (see §1.2 and §3). The system \tilde{F}_{-u_n} is nonsingular over the subfaces under consideration. In the actual implementation, LINPACK routines DGECC and DGESL were used to compute the preconditioners.
4. In the search phase, in Step 1f (or 2f), we first expand the subinterval \mathbf{y}_n^0 (or \mathbf{x}_n^0) by $\epsilon = \frac{1}{2}w(\mathbf{y}_n^0)$ at both ends. If the existence and uniqueness of the zero of \tilde{F}_{-u_n} can be verified over the corresponding subface, then we expand the subinterval by 2ϵ at both ends, and then 4ϵ and so on until the existence and uniqueness of the zero can not be established.
5. In the search phase, in Step 1g (or 2g), the underlying idea is that the farther away the interval \mathbf{y}_n^l is from the interval \mathbf{y}_n^0 whose corresponding subface of \mathbf{x}_n (or $\mathbf{x}_{\bar{n}}$) contains a unique solution of $\tilde{F}_{-u_n} = 0$ and/or the narrower the interval \mathbf{y}_n^l is, the more probable it is that we can verify that $\tilde{F}_{-u_n} \neq 0$ over the subface of \mathbf{x}_n (or $\mathbf{x}_{\bar{n}}$) which corresponds to \mathbf{y}_n^l .

6.2. Computational Complexity.

Derivation of the Computational Complexity

Box-setting Phase: Step 1 is of order $\mathcal{O}(n^3)$. Step 2 is of order $\mathcal{O}(n)$. Step 3 is of order $\mathcal{O}(n^2)$.

Thus, the order of this phase is $\mathcal{O}(n^3)$.

Elimination Phase: Step 1ai and 1bi are of order $\mathcal{O}(n^2)$. Step 1aii and 1bii are of order $\mathcal{O}(1)$.

Thus, the order of this phase is $\mathcal{O}(n^3)$.

Search Phase: Step 1a and 2a are of order $\mathcal{O}(n^3)$. Step 1b and 2b are of order $\mathcal{O}(n^2)$. Step 1c and 2c are of order $\mathcal{O}(1)$. Step 1d and 2d are of order $\mathcal{O}(n^3)$. Step 1e and 2e are of order $\mathcal{O}(n^3)$. Step 1f and 2f are of order $N_{infl} * \mathcal{O}(n^3)$. (See explanation below.) Step 1g and 2g are of order $N_{proc} * \mathcal{O}(n^3)$. (See explanation below.) Step 1h and 2h are of order $\mathcal{O}(n^2)$. Step 1i and 2i are of order $\mathcal{O}(n^3)$. Step 1j and 2j are of order $\mathcal{O}(1)$. The last step of this phase is of order $\mathcal{O}(1)$ too.

Thus, the order of this phase is $\mathcal{O}(n^3)$.

The order of the overall algorithm is thus $\mathcal{O}(n^3)$.

Remarks

1. There are two performance measures in the algorithm, N_{infl} and N_{proc} . N_{infl} is the number of inflations the algorithm did in Step 1f or 2f. N_{proc} is the number of subintervals of $\mathbf{y}_n \setminus \mathbf{y}_n^1$ the algorithm processed in Step 1g or subintervals of $\mathbf{x}_n \setminus \mathbf{x}_n^1$ the algorithm processed in Step 2g, i.e. number of \mathbf{y}_n^l 's plus number of \mathbf{y}_n^u 's in Step 1g or number of \mathbf{x}_n^l 's plus number of \mathbf{x}_n^u 's in Step 2g. (See the algorithm in §6.1).
2. The order of the algorithm cannot be improved, since computing preconditioners of the original system and the system \tilde{F}_{-u_n} is necessary and computing each preconditioner is of order $\mathcal{O}(n^3)$.

7. Test Problems and Test Environment.

7.1. Test Problems. Before describing the test set, we introduce one more problem. Motivated by [6, Lemma 2.4], we considered systems of the following form.

EXAMPLE 4. *Set*

$$f(x) = h(x, t) = (1 - t)(Ax - x^2) - tx,$$

where $A \in \mathbb{R}^{n \times n}$ is the matrix corresponding to central difference discretization of the boundary value problem $-u'' = 0$, $u(0) = u(1) = 0$ and $x^2 = (x_1^2, \dots, x_n^2)^T$. t was chosen to be equal to

$$t_1 = \lambda_1 / (1 + \lambda_1),$$

where λ_1 is the largest eigenvalue of A .

The homotopy h in Example 4 has a simple bifurcation point at $t = t_1$, where the two paths cross obliquely. That is, there are two solutions to $f(x) = 0$ near $x = 0$, for all t near t_1 and on either side of t_1 . Furthermore, the quadratic terms in the Taylor expansion for f do not vanish at $t = t_1$.

The test set consists of Example 2, Example 3 with $\epsilon = +10^{-6}$ and -10^{-6} , and Example 4 with $n = 5, 10, 20, 40, 80, 160, 320$. For all the test problems, we used $(0, 0, \dots, 0)$ as a good approximate solution to the problem $F(x) = 0$. Actually, it's the exact solution in Example 2 and Example 4. $w(\mathbf{x}_k)$ and $w(\mathbf{y}_k)$ were set to 10^{-3} for $1 \leq k \leq n - 1$. $w(\mathbf{x}_n)$ and $w(\mathbf{y}_n)$ were computed automatically by the algorithm. In fact, $w(\mathbf{x}_k)$ and $w(\mathbf{y}_k)$, $1 \leq k \leq n - 1$, can also be computed automatically by the algorithm, depending on the accuracy of the approximate solution. At present, we used the known true solutions to Example 2 and Example 4 and the known approximate solution to Example 3 to test the algorithm and set the widths apparently small but otherwise arbitrary.

For all the problems, the algorithm succeeded and returned a degree of 2.

7.2. Test Environment. The algorithm in §6.1 was programmed in the Fortran 90 environment developed and described in [8, 9]. Similarly, all the functions in the test problems were programmed using the same Fortran 90 system, and internal symbolic representations of the functions were generated prior to execution of the numerical tests. In the actual tests, generic routines then interpreted the internal representations to obtain both floating point and internal values.

The LINPACK routines DGECCO and DGESL were used in Step 1 of the box-setting phase, and in Step 1e, 2e, 1f and 2f of the search phase to compute the preconditioners. (See the algorithm and its explanation in §6.1.)

The Sun Fortran 90 compiler version 1.2 was used on a Sparc Ultra model 140 with optimization level 0. Execution times were measured using the routine DSECND. All times are given in CPU seconds.

8. Numerical Results. We present the numerical results in Table 8.1 and some statistical data in Table 8.2.

The column labels of Table 8.1 are as follows.

Problem: names of the problems identified in §7.1.

n : number of independent variables.

Success: whether the algorithm was successful.

Degree: topological degree returned by the algorithm.

CPU Time: CPU time in seconds of the algorithm.

TABLE 8.1
Numerical Results

Problem	n	Success	Degree	CPU Time	Time Ratio
Example 2	2	Yes	2	0.0761	
Example 3 ($\epsilon = +10^{-6}$)	2	Yes	2	0.0511	
Example 3 ($\epsilon = -10^{-6}$)	2	Yes	2	0.0513	
Example 4	5	Yes	2	0.6806	
Example 4	10	Yes	2	3.3403	4.91
Example 4	20	Yes	2	19.440	5.82
Example 4	40	Yes	2	140.34	7.22
Example 4	80	Yes	2	1123.6	8.01
Example 4	160	Yes	2	8891.3	7.91
Example 4	320	Yes	2	65395.5	7.36

Time Ratio: This only applies to Example 4. It's the ratio of two successive CPU times.

The column labels of Table 8.2 are as follows.

Problem: names of the problems identified in §7.1.

n : number of independent variables.

N_{infl} : number of inflations the algorithm did in Step 1f or 2f for the indicated face \mathbf{x}_n , $\mathbf{x}_{\bar{n}}$, \mathbf{y}_n or $\mathbf{x}_{\bar{y}}$ (see explanation of the computational complexity of the algorithm in §6.2).

N_{proc} : number of intervals the algorithm processed in Step 1g or 2g for the indicated face \mathbf{x}_n , $\mathbf{x}_{\bar{n}}$, \mathbf{y}_n or $\mathbf{x}_{\bar{y}}$ (see explanation of the computational complexity of the algorithm in §6.2).

We can see from Table 8.1 that the algorithm was successful on all the problems in the test set. The algorithm is of order $\mathcal{O}(n^3)$. But there are many $\mathcal{O}(n^3)$ and $\mathcal{O}(n^2)$ steps. Some steps have many $\mathcal{O}(n^3)$ and $\mathcal{O}(n^2)$ substeps, and some of the substeps still have many $\mathcal{O}(n^2)$ structures. So, when n was small, those lower order structures had significant influence on the CPU time. But, when n was large, the $\mathcal{O}(n^3)$ terms dominated. We can see this from the time ratios of Example 4 in Table 8.1.

In Table 8.2, for all the problems, there were two faces of \mathbf{x}_n , $\mathbf{x}_{\bar{n}}$, \mathbf{y}_n and $\mathbf{y}_{\bar{n}}$ for which $N_{infl} = 0$. This is because the algorithm verified that $u_n < 0$ on each of those two faces in Step 1b or 2b, and returned a degree contribution of each of those two faces as 0. So, the algorithm didn't proceed to Step 1f or 2f. For the same reason, $N_{proc} = 0$ for those two faces. For the remaining two faces for which the algorithm did proceed to Step 1f or 2f, N_{infl} is small.

In Step 1g or 2g which immediately follows the inflations, $N_{proc} = 0$ for Example 2 and Example 3. This is because the inflations had covered the whole interval \mathbf{y}_n . (See the algorithm in §6.1 and the explanation of its computational complexity in §6.2.) More significant is that $N_{proc} = 2$ in Example 4 regardless of small n or large n . This is because only one interval was processed to verify that $\tilde{F}_{-u_n} = 0$ has no solutions when $x_n \in \mathbf{x}_n^l$ and only one interval was processed to verify that $\tilde{F}_{-u_n} = 0$ has no solutions when $x_n \in \mathbf{x}_n^u$. (See the algorithm in §6.1 and the explanation of its computational complexity in §6.2.) This means that the algorithm was quite efficient.

TABLE 8.2
Statistical Data

Problem	n	N_{infl}				N_{proc}			
		\underline{x}_n	\overline{x}_n	\underline{y}_n	\overline{y}_n	\underline{x}_n	\overline{x}_n	\underline{y}_n	\overline{y}_n
Example 2	2	6	6	0	0	0	0	0	0
Example 3 ($\epsilon = +10^{-6}$)	2	2	2	0	0	0	0	0	0
Example 3 ($\epsilon = -10^{-6}$)	2	2	2	0	0	0	0	0	0
Example 4	5	0	0	5	5	0	0	2	2
Example 4	10	0	0	5	5	0	0	2	2
Example 4	20	0	0	4	4	0	0	2	2
Example 4	40	0	0	4	4	0	0	2	2
Example 4	80	0	0	4	4	0	0	2	2
Example 4	160	0	0	4	4	0	0	2	2
Example 4	320	0	0	3	3	0	0	2	2

9. Conclusions and Future Work. When we tested the algorithm, we took advantage of knowing the true solutions. (See §7.1.) For this reason, we set $w(\mathbf{x}_k)$ and $w(\mathbf{y}_k)$, $1 \leq k \leq n-1$ somewhat arbitrarily. But we plan to have the algorithm eventually compute these, based on the accuracy of the approximate solution obtained by a floating point algorithm and the accuracy of the quadratic model.

We presented an algorithm which was designed to work for the case that the rank deficiency of the Jacobian matrix at the singular solution is one. But the analysis in §5 and the algorithm in §6.1 can be generalized to general rank deficiency. Also, at present, it is assumed that the second derivatives $\frac{\partial^2 f_n}{\partial x_k \partial x_l}$, $1 \leq k \leq n, 1 \leq l \leq n$ don't vanish simultaneously at the singular solution. In fact, the analysis in §5 and the algorithm in §6.1 can be generalized to the general case that the derivatives of f_n of order 1 through r ($r \geq 2$) vanish simultaneously at the singular solution. But computing higher order derivatives may be expensive. Those two generalizations can also be combined, i.e. any rank deficiency and any order of derivatives of f_n that vanish. We will pursue these generalizations in the future.

Another future direction of this study is to apply the algorithms to bifurcation problems and other physical models.

REFERENCES

- [1] O. ABERTH, *Computation of topological degree using interval arithmetic, and applications*, Math. Comp., 62 (1994), pp. 171–178.
- [2] G. ALEFELD AND J. HERZBERGER, *Introduction to Interval Computations*, Academic Press, New York, 1983.
- [3] P. ALEXANDROFF AND H. HOPF, *Topologie*, Chelsea, 1935.
- [4] J. CRONIN, *Fixed Points and Topological Degree in Nonlinear Analysis*, American Mathematical Society, Providence, RI, 1964.
- [5] E. R. HANSEN, *Global Optimization Using Interval Analysis*, Marcel Dekker, Inc., New York, 1992.
- [6] H. JÜRGENS, H.-O. PEITGEN, AND D. SAUPE, *Topological perturbations in the numerical nonlinear eigenvalue and bifurcation problems*, in Analysis and Computation of Fixed Points, S. M. Robinson, ed., New York, 1980, Academic Press, pp. 139–181.

- [7] R. B. KEARFOTT, *Computing the Degree of Maps and a Generalized Method of Bisection*, PhD thesis, University of Utah, Department of Mathematics, 1977.
- [8] ———, *A Fortran 90 environment for research and prototyping of enclosure algorithms for non-linear equations and global optimization*, ACM Trans. Math. Software, 21 (1995), pp. 63–78.
- [9] ———, *Rigorous Global Search: Continuous Problems*, Kluwer, Dordrecht, Netherlands, 1996.
- [10] A. NEUMAIER, *Interval Methods for Systems of Equations*, Cambridge University Press, Cambridge, England, 1990.
- [11] H. RATSCHKE AND J. ROKNE, *New Computer Methods for Global Optimization*, Wiley, New York, 1988.
- [12] S. M. RUMP, *Verification methods for dense and sparse systems of equations*, in Topics in Validated Computations, J. Herzberger, ed., Amsterdam, 1994, Elsevier Science Publishers, pp. 63–135.
- [13] F. STENGER, *Computing the topological degree of a mapping in \mathbb{R}^n* , Numer. Math., 25 (1975), pp. 23–38.
- [14] M. STYNES, *An Algorithm for the Numerical Calculation of the Degree of a Mapping*, PhD thesis, Oregon State University, Department of Mathematics, Corvallis, Oregon, 1977.