



ELSEVIER

Applied Numerical Mathematics 13 (1994) 453–468

APPLIED
NUMERICAL
MATHEMATICS

An improvement of Gargantini's simultaneous inclusion method for polynomial roots by Schröder's correction

Carsten Carstensen^{a,*}, Miodrag S. Petković^b

^a Department of Mathematics, Heriot-Watt University, Edinburgh EH14 4AS, UK

^b Faculty of Electronic Engineering, Department of Mathematics, University of Niš, Beogradska 14, P.O. Box 73, 18000 Niš, Yugoslavia

Abstract

The interval version of the complex third-order method of Maehly, Börsch-Supan, Ehrlich, and Aberth is the most efficient method for simultaneous inclusion of simple polynomial roots [14]. In this note, Gargantini's generalization of this third-order interval method for multiple roots is accelerated using Schröder's modification of Newton's corrections and modifying the required interval inversions. The underlying idea is that the iteration of the midpoints of the interval method should be similar to Nourin's acceleration of the above-mentioned complex third-order method improving the convergence of the midpoints. Since the convergence of the radii and the midpoints are coupled it can be proved that the R-orders of convergence of the radii of the newly presented Schröder-like interval methods are asymptotically greater than 3.5. Hence two of these methods are more efficient than the most efficient one known before. Numerical results and an analysis of computational efficiency are included.

Key words: Polynomial zeros; Simultaneous root finding method; Interval methods

1. Introduction

During the last twenty years various authors developed the techniques for *a posteriori* error estimates for the approximation of polynomial zeros. These devices mostly use Gerschgorin's theorem, Rouché's theorem, the fixed point principle, and so on. One of the simplest methods for the estimate of upper error bounds of the produced approximations z_1, \dots, z_n to the zeros of a polynomial P consists of the combination of a suitable iterative method realized in ordinary complex arithmetic and the so-called inclusion disc of the form $|z_j - z| \leq r(z_1, \dots, z_n)$ which contains at least one zero of P (see [9] or [14, Chapter 6]). A quite different approach to error estimates for a given set of approximate zeros uses circular arithmetic, as pointed out by

* Corresponding author.

Gargantini and Henrici [8]. Simultaneous iterative methods, realized in circular arithmetic, yield resulting discs containing the complex zeros in *each* iterative step. In this manner the upper bounds for the zeros, given by the radii of discs, are obtained automatically. A further advantage of inclusion methods lies in the possibility of taking into account rounding errors by applying *rounding interval arithmetic* which is built in modern high-level programming languages PASCAL-XSC [10] and ACRITH-XSC [2]. It should be noted that nowadays the computational costs of this type of arithmetic is only slightly greater compared with standard floating-point arithmetic. Therefore, a reasonably high computational efficiency of inclusion methods, together with the very useful property of self-verifying results, makes these methods to be often implemented in many problems of applied mathematics and techniques.

Let P be a monic polynomial of degree $n \geq 3$ with distinct zeros ζ_1, \dots, ζ_m ($m \leq n$) of multiplicities $\mu_1, \dots, \mu_m \geq 1$ ($\mu_1 + \dots + \mu_m = n$), that is

$$P(z) = z^n + a_1 z^{n-1} + \dots + a_{n-1} z + a_n = (z - \zeta_1)^{\mu_1} (z - \zeta_2)^{\mu_2} \dots (z - \zeta_m)^{\mu_m}.$$

Applying the logarithmic derivatives we find

$$\frac{P'(z)}{P(z)} = \sum_{k=1}^m \mu_k (z - \zeta_k)^{-1},$$

wherefrom

$$\zeta_j = z - \frac{1}{\frac{1}{N_j} - \frac{1}{\mu_j} \sum_{k=1, k \neq j}^m \mu_k (z - \zeta_k)^{-1}}, \quad j = 1, \dots, m. \tag{1}$$

Here $N_j := \mu_j P(z) / P'(z)$ is Schröder's modification of Newton's correction for multiple roots.

Assume that $m \geq 2$ disjoint discs Z_1, \dots, Z_m in the complex plane are found such that $\zeta_j \in Z_j$ for any $j \in \{1, \dots, m\}$. Starting from (1) and using circular interval operations (see [3, Chapters 5 and 6]) Gargantini established in [7] the Schröder-like algorithm for the simultaneous inclusion of all zeros of P ,

$$\hat{Z}_j = z_j - \frac{1}{\frac{1}{N_j} - \frac{1}{\mu_j} \sum_{k=1, k \neq j}^m \mu_k (z_j - Z_k)^{-1}}, \quad j = 1, \dots, m. \tag{2}$$

The inclusion method (2) has a *cubic* convergence. An algorithm of the form (2) in ordinary complex arithmetic for the case of simple zeros was considered before by Maehly [11], Börsch-Supan [4], Ehrlich [6], and Aberth [1]. Using Newton's correction Nourein [12] improved this algorithm and stated the iterative formula

$$\hat{z}_j = z_j - \frac{1}{\frac{1}{N_j} - \frac{1}{\mu_j} \sum_{k=1, k \neq j}^m \mu_k (z_j - z_k + N_k)^{-1}}, \quad j = 1, \dots, m, \tag{3}$$

where the order of convergence equals *four*.

The aim of this paper is to improve the third-order interval method (2) following Nourain's idea, that is, using Newton's correction. The convergence of the midpoints and the convergence of the radii are coupled (cf. Section 4 below) such that a further improvement of the midpoints via (3) actually improves the convergence of the radii. To ensure that the midpoints are determined through (3) we have to enlarge the interval inversions. A careful analysis shows that the enlarged radii do not disturb the mentioned improvement. We will use three types of inversion of a disc Z to construct various inclusion algorithms. The new algorithms have the R-order of convergence of the radii greater than 3.5 depending on the applied inversion (for the definition of the R-order see [13,14]). They do not require additional calculations because already found values of N_j are used. Consequently, these algorithms possess a great computational efficiency (Section 6): two of them are even more efficient than the interval method (2) which used to be the most efficient known method for the simultaneous inclusion of polynomial zeros [14, Chapter 6].

Using the notation $Z = \{z, r\}$, where $z = \text{mid}(Z)$ and $r = \text{rad}(Z)$ are the center and radius of Z , we consider the discs (assuming that $0 \notin Z$, that is, $|z| > r$):

$$Z^{-1} := \{z, r\}^{-1} := \left\{ \frac{1}{z \cdot \left(1 - \frac{r^2}{|z|^2}\right)}, \frac{r}{|z|^2 - r^2} \right\},$$

$$Z^{I_1} := \{z, r\}^{I_1} := \left\{ \frac{1}{z}, \frac{r}{|z| \cdot (|z| - r)} \right\},$$

$$Z^{I_2} := \{z, r\}^{I_2} := \left\{ \frac{1}{z}, \frac{2r}{|z|^2 - r^2} \right\}.$$

Sometimes, we will write $\{z, r\}^I$ instead of $\{z, r\}^{-1}$ (Sections 5 and 6). It is not hard to see that $Z^{-1} \subseteq Z^{I_1} \subseteq Z^{I_2}$. Among the above inversions only $\{z, r\}^{-1}$ is the exact operation, that is, $\{z, r\}^{-1} = \{z^{-1}: z \in Z\}$, but in general $\text{mid}(Z^{-1}) \neq \text{mid}(Z)^{-1}$. If $0 \notin Z = \{z, r\}$ and $\text{INV}_i(Z)$ denotes one of the three inversions Z^{-1} , Z^{I_1} , Z^{I_2} , then there holds the inclusion $Z^{-1} \subseteq \text{INV}_i(Z)$ and the estimates

$$|\text{mid}(\text{INV}_i(Z))| \leq \frac{|z|}{|z|^2 - r^2}, \quad \text{rad}(\text{INV}_i(Z)) \leq \frac{2r}{|z|^2 - r^2}. \quad (4)$$

We emphasize that the three cases $\text{INV}_i = ()^{-1}$, $\text{INV}_i = ()^{I_1}$, and $\text{INV}_i = ()^{I_2}$, respectively, are handled simultaneously and only the properties (4) of INV_i will be used below if not stated explicitly otherwise.

2. The methods

To compute circular approximations for the distinct zeros ζ_1, \dots, ζ_m of a polynomial P simultaneously, we assume that $m \geq 2$ disjoint discs Z_1, \dots, Z_m containing the zeros ζ_1, \dots, ζ_m are known.

Writing $\text{mid}(Z_j) =: z_j$ and $\text{rad}(Z_j) =: r_j$ for the center and the radius of the disc Z_j , one step of the new Schröder-like algorithms with Nourein's approach reads $(Z_1, \dots, Z_m) \mapsto (\hat{Z}_1, \dots, \hat{Z}_m)$ with

$$\hat{Z}_j := z_j - \text{INV}_1 \left(\frac{1}{N_j} - \frac{1}{\mu_j} \sum_{k=1, k \neq j}^m \mu_k \cdot \text{INV}_2(z_j - Z_k + N_k) \right), \quad j = 1, \dots, n, \quad (5)$$

where the complex number $N_j := \mu_j P(z_j) / P'(z_j)$ is Schröder's modification of Newton's correction for multiple roots for the center of Z_j .

In (5) INV_1 and INV_2 denote inversions of a disc defined in Section 1, that is, $\text{INV}_1 \in \{()^{-1}, ()^{I_1}, ()^{I_2}\}$ and $\text{INV}_2 \in \{()^{-1}, ()^{I_1}, ()^{I_2}\}$; thus (5) describes six different methods simultaneously.

The convergence analysis of the new algorithms is presented in Sections 3 and 4, while practical aspects (numerical results and computational efficiency) are the subject of Sections 5 and 6.

3. Convergence results

The following theorems state that method (5) is a locally convergent interval method having the R-order of convergence $\frac{1}{2}(3 + \sqrt{17}) \approx 3.562$ or 4, depending on the choice of $\text{INV}_i \in \{()^{-1}, ()^{I_2}, ()^{I_2}\}$. The proofs are given in the next section.

Theorem 3.1. *Let $(Z_1, \dots, Z_n) =: (Z_1^{(0)}, \dots, Z_n^{(0)})$ be initial discs such that $\zeta_j \in Z_j$ ($j = 1, \dots, n$) and*

$$\frac{r}{d} \leq \frac{1}{4n}, \quad (6)$$

where

$$r := \max_{j=1, \dots, m} \text{rad}(Z_j),$$

$$d := \min_{i, j=1, \dots, m, i \neq j} |\text{mid}(Z_i) - \text{mid}(Z_j)|.$$

Then, the method (5) is feasible, i.e. it defines a sequence of discs $(Z_j^{(\nu)} \mid j = 1, \dots, m)_{\nu=0,1,2,\dots}$, for any $j \in \{1, \dots, m\}$ and $\nu = 0, 1, 2, \dots$, there holds

$$\zeta_j \in Z_j^{(\nu)},$$

and the sequence of radii $(\text{rad}(Z_j^{(\nu)}))_{\nu=0,1,2,\dots}$ tends towards zero.

Theorem 3.2. *Let $O((5))$ denote the R-order of convergence of the radii for method (5), where $\text{INV}_1, \text{INV}_2 \in \{()^{-1}, ()^{I_1}, ()^{I_2}\}$. Then,*

$$O((5)) \geq \begin{cases} \frac{1}{2}(3 + \sqrt{17}), & \text{if } \text{INV}_2 = ()^{-1}, \\ 4, & \text{otherwise.} \end{cases}$$

Theorem 3.2 shows the improvement of the R-order of convergence from 3 of (2) to 3.562 or 4 of method (5) whereas the computer effort is comparable (cf. Section 6). It seems to be surprising that an enlargement of the inversions in (5) gives better convergence. It is caused by the better convergence of the midpoints yielding smaller corrections N_j for instance and hence smaller radii, cf. the proof of Theorem 3.2 below dealing with the coupling of the convergence of the midpoints and the radii.

4. Proofs of the convergence theorems

The proofs of Theorem 3.1 and 3.2 are divided in several lemmas.

The proof of Theorem 3.1 is by induction on ν and we consider the typical step for $\nu = 0$ first (neglecting the iteration index ν). In addition to the notation of the theorems, introduce some abbreviations. For any $j, k \in \{1, \dots, m\}$, let

$$z_j := \text{mid}(Z_j),$$

$$r_j := \text{rad}(Z_j) \leq r,$$

$$\varepsilon_j := z_j - \zeta_j,$$

$$a_j := \frac{1}{\mu_j} \sum_{k=1, k \neq j}^m \frac{\mu_k}{z_j - \zeta_k},$$

$$N_j := \frac{\mu_j P(z_j)}{P'(z_j)} = (1/\varepsilon_j + a_j)^{-1} \left(\text{since } \frac{P'(z)}{P(z)} = \sum_{k=1}^m \mu_k (z - \zeta_k)^{-1} \right),$$

$$v_{jk} := z_j - z_k + N_k.$$

If $\hat{Z}_1, \dots, \hat{Z}_m$ are given by (5) having the centers $\hat{z}_1, \dots, \hat{z}_m$ and the radii $\hat{r}_1, \dots, \hat{r}_m$, let

$$\hat{r} := \max_{j=1, \dots, m} \hat{r}_j, \quad \hat{d} := \min_{j, k=1, \dots, m, j \neq k} |\hat{z}_j - \hat{z}_k|.$$

Finally, for $j \in \{1, \dots, m\}$, let

$$W_j := \frac{1}{N_j} - \frac{1}{\mu_j} \sum_{k=1, k \neq j}^m \mu_k \cdot \text{INV}_2(z_j - z_k + N_k) =: \{u_j, \rho_j\}.$$

Lemma 4.1. *If $d/r \geq 4n$ and $\zeta_j \in Z_j$ for all $j \in \{1, \dots, m\}$, then the inversions in (5) exist (i.e. $0 \notin W_j$, $0 \notin z_j - z_k + N_k$, for $j, k \in \{1, \dots, m\}$, $j \neq k$) and, for any $j \in \{1, \dots, m\}$, there holds*

$$\rho_j \leq \frac{2(n-1)r}{(d-3r)(d-r)}, \tag{7a}$$

$$|u_j| \geq \frac{3}{4|\varepsilon_j|} - \frac{(n-1)(d-2r)}{(d-3r)(d-r)}, \tag{7b}$$

$$\zeta_j \in \hat{Z}_j, \tag{7c}$$

$$\hat{r} \leq \frac{r}{n-1}, \tag{7d}$$

$$\hat{d} \geq d - 5r, \tag{7e}$$

$$\frac{\hat{r}}{\hat{d}} \leq \frac{6r}{7d}, \tag{7f}$$

$$\hat{Z}_1, \dots, \hat{Z}_m \text{ are pairwise disjoint.} \tag{7g}$$

Proof. Note that $\zeta_j \in Z_j$ yields $|\varepsilon_j| \leq r_j \leq r$. Thus,

$$\begin{aligned} |a_j| &\leq \sum_{k=1, k \neq j}^m \frac{\mu_k}{|z_j - z_k| - |z_k - \zeta_k|} \\ &\leq \frac{n-1}{d-r} = \frac{n-1}{r(d/r-1)} \\ &\leq \frac{n-1}{r(4n-1)} < \frac{1}{4r}, \end{aligned}$$

which implies

$$|a_j \varepsilon_j| < \frac{1}{4}.$$

Therefore, using $N_j^{-1} = 1/\varepsilon_j + a_j$,

$$\begin{aligned} |z_j - N_j - \zeta_j| &= |\varepsilon_j - N_j| = \left| \varepsilon_j - \frac{\varepsilon_j}{1 + \varepsilon_j a_j} \right| \\ &\leq |\varepsilon_j|^2 \frac{|a_j|}{1 - |\varepsilon_j a_j|} \leq \frac{4}{3} |\varepsilon_j|^2 \cdot |a_j| < \frac{1}{3} |\varepsilon_j| < r, \end{aligned}$$

whence $\zeta_j \in Z_j - N_j$. Note also that

$$|N_j| \leq |z_j - N_j - \zeta_j| + |z_j - \zeta_j| \leq \frac{4}{3} |\varepsilon_j| \leq \frac{4}{3} r.$$

Hence, for all $j, k \in \{1, \dots, m\}$, $j \neq k$,

$$|v_{jk}| \geq |z_j - z_k| - |N_k| \geq d - \frac{4}{3} r > d - 2r.$$

Since $d \geq 4nr$, we still conclude $|v_{jk}| > d - 2r > r \geq r_k$, whence $0 \notin \{v_{jk}, r_k\} = z_j - Z_k + N_k$.

Therefore, W_j exists and we can estimate $u_j := \text{mid}(W_j)$ and $\rho_j := \text{rad}(W_j)$ in the sequel.

Proof of (7a). By (4),

$$\begin{aligned} \rho_j &= \frac{1}{\mu_j} \sum_{k=1, k \neq j}^m \mu_k \cdot \text{rad}(\text{INV}_2(\{v_{jk}, r_k\})) \\ &\leq \frac{1}{\mu_j} \sum_{k=1, k \neq j}^m \mu_k \cdot \frac{2r}{|v_{jk}|^2 - r^2} \\ &\leq (n-1) \frac{2r}{(d-2r)^2 - r^2} = \frac{2r(n-1)}{(d-r)(d-3r)}. \end{aligned}$$

Proof of (7b). Again, by (4),

$$\begin{aligned} |u_j| &\geq \frac{1}{|N_j|} - \frac{1}{\mu_j} \sum_{k=1, k \neq j}^m \mu_k \cdot |\text{mid}(\text{INV}_2(z_j - Z_k + N_k))| \\ &\geq \frac{3}{4|\varepsilon_j|} - \sum_{k=1, k \neq j}^m \mu_k \cdot \frac{|v_{jk}|}{|v_{jk}|^2 - r_k^2} \\ &\geq \frac{3}{4|\varepsilon_j|} - (n-1) \frac{d-2r}{(d-2r)^2 - r^2} \\ &= \frac{3}{4|\varepsilon_j|} - \frac{(n-1)(d-2r)}{(d-r)(d-3r)}. \end{aligned}$$

We continue proving that \hat{Z}_j exists. By (7a) and (7b), for $j \in \{1, \dots, m\}$,

$$\begin{aligned} |u_j|^2 - \rho_j^2 &\geq \left(\frac{3}{4r} - \frac{(n-1)(d-2r)}{(d-r)(d-3r)} \right)^2 - \left(\frac{2r(n-1)}{(d-r)(d-3r)} \right)^2 \\ &= \frac{((3/4r)(d-r)(d-3r) - (n-1)(d-2r))^2 - (2r(n-1))^2}{(d-r)^2(d-3r)^2} \\ &\geq \frac{(\frac{3}{4}(d-r)(4n-3) - (n-1)(d-2r))^2 - (2r(n-1))^2}{(d-r)^2(d-3r)^2} \\ &\geq (n-1)^2 \frac{(3(d-r) - (d-2r))^2 - 4r^2}{(d-r)^2(d-3r)^2} \\ &\geq (n-1)^2 \frac{(2d-4r)^2 - 4r^2}{(d-r)^2(d-3r)^2} \\ &= \frac{4(n-1)^2}{(d-r)(d-3r)}. \end{aligned}$$

Consequently, since $d/r \geq 4n \geq 12$, $|u_j|^2 - \rho_j^2 > 0$, whence $0 \notin W_j$. Altogether, the discs which must be inverted in (5) do not contain zero; thus method (5) is feasible.

Proof of (7c). Since $\zeta_k \in Z_k - N_k$ for any $k \in \{1, \dots, n\}$ by the inclusion property from (1) we obtain $\zeta_j \in \hat{Z}_j$, i.e. (7c).

Proof of (7d). By (4), (7a), and the above estimate of $|u_j|^2 - \rho_j^2$ there holds for any $j \in \{1, \dots, m\}$

$$\hat{r}_j \leq \frac{2\rho_j}{|u_j|^2 - \rho_j^2} \leq 2 \frac{2r(n-1)}{4(n-1)^2} = \frac{r}{n-1}.$$

Proof of (7e). By (7a) and (7b), for any $j \in \{1, \dots, m\}$,

$$\begin{aligned} \frac{\rho_j}{|u_j|} &\leq \frac{2r(n-1)}{(3/4r)(d-r)(d-3r) - (n-1)(d-2r)} \\ &\leq \frac{2r}{3(d-r) - (d-2r)} = \frac{r}{d - \frac{1}{2}r} \\ &\leq \frac{1}{4(n-1)}. \end{aligned}$$

Therefore, and by (4) and (7b)

$$\begin{aligned} |\hat{z}_j - z_j| &= |\text{mid}(\text{INV}_1(\{u_j, \rho_j\}))| \leq \frac{1}{|u_j|(1 - \rho_j^2/|u_j|^2)} \\ &\leq \frac{r}{\left(\frac{3}{4} - \frac{(n-1)(d-2r)r}{(d-r)(d-3r)}\right) \cdot \left(1 - \frac{1}{16(n-1)^2}\right)} \\ &\leq \frac{r}{\frac{1}{2} \cdot \left(1 - \frac{1}{16(n-1)^2}\right)} \leq \frac{2r}{1 - \frac{1}{64}} = \frac{128}{63}r. \end{aligned}$$

Thus, for some particular $j, k \in \{1, \dots, m\}$, $j \neq k$

$$\begin{aligned} \hat{d} &= |\hat{z}_j - \hat{z}_k| \\ &\geq |z_j - z_k| - |\hat{z}_j - z_j| - |\hat{z}_k - z_k| \geq d - 5r \end{aligned}$$

Proof of (7f). By (7d) and (7e)

$$\begin{aligned} \frac{\hat{r}}{\hat{d}} &\leq \frac{r/(n-1)}{d-5r} \leq \frac{r}{d(n-1)(1-5/4n)} \\ &= \frac{r}{d} \frac{n}{(n-1)(n-\frac{5}{4})} \leq \frac{6}{7} \cdot \frac{r}{d} < \frac{r}{d} < \frac{1}{4n} \quad \text{for } n \geq 3. \end{aligned}$$

Proof of (7g). We mention that $d/r \geq 4n$ implies that Z_1, \dots, Z_m are pairwise disjoint which immediately follows from $d - 2r > 0$. \square

Proof of Theorem 3.1. From Lemma 4.1 we conclude by induction on $\nu = 0, 1, 2, \dots$ that method (5) defines discs $Z_j^{(\nu)}$ including ζ_j for any $\nu = 0, 1, 2, \dots$, $j \in \{1, \dots, m\}$, such that

$$\frac{\max_{j=1, \dots, m} \text{rad}(Z_j^{(\nu)})}{\min_{j, k=1, \dots, m, j \neq k} |\text{mid}(Z_j^{(\nu)}) - \text{mid}(Z_k^{(\nu)})|} \leq \frac{1}{4n}.$$

Therefore the conclusions of Lemma 4.1 hold in any iteration step. In particular, (7d) shows that method (5) converges. \square

Let $(Z_j^{(\nu)} | j = 1, \dots, m)_{\nu=0,1,2,\dots}$ denote the sequence of discs generated by method (5). Assume that the radii converge towards zero such that for large ν there holds (6). Without loss of generality, let us assume that (6) holds for $\nu = 0$ (and hence for all ν). Moreover, it will be used several times that, going to determine the R-order of convergence, we may assume that all the discs under consideration are sufficiently small (otherwise consider the sequence related to the indices $\nu = \nu_0, \nu_0 + 1, \nu_0 + 2, \dots$ having the same R-order of convergence where ν_0 is sufficiently large).

Since the zeros ζ_1, \dots, ζ_m are fixed and included in the discs $Z_1^{(\nu)}, \dots, Z_m^{(\nu)}$ various denominators in Lemma 4.1 are bounded. Hence we may use the Landau symbols $O()$ to suppress the bounds and stress the asymptotical behavior.

For expressions $\text{term1}(j, \nu)$ and $\text{term2}(j, \nu)$ (depending on j, ν, P , and the initial discs) we define

$$\text{term1}(j, \nu) = O(\text{term2}(j, \nu))$$

$$\text{iff } \max_{j=1,\dots,m} \sup_{\nu=0,1,2,\dots} \left| \frac{\text{term1}(j, \nu)}{\text{term2}(j, \nu)} \right| < \infty.$$

For example, using the notation from Lemma 4.1 but adding the iteration index ν , (7a) and (7b) read

$$\rho_j^{(\nu)} = O(r^{(\nu)}), \quad |\varepsilon_j^{(\nu)} u_j^{(\nu)}| = O(1),$$

respectively.

In the above analysis we used the rough estimate $|\varepsilon_j| \leq r_j \leq r$ but we controlled the constants to prove that (6) guarantees convergence. Conversely, we will neglect the constants but discuss in detail the dependence on ε_j and r_j . Therefore, define for any $\nu = 0, 1, 2, \dots$

$$E_\nu := \max_{j=1,\dots,m} |\varepsilon_j^{(\nu)}|, \quad R_\nu := \max_{j=1,\dots,m} |r_j^{(\nu)}| = r^{(\nu)}.$$

If we neglect the iteration index and write E instead of E_ν , then \hat{E} denotes $E_{\nu+1}$.

Lemma 4.2. $R_{\nu+1} = O(R_\nu E_\nu^2)$.

Proof. Neglecting the iteration index ν we see from (7a) and (7b) that

$$\rho = O(r), \quad |u_j| = O(1/|\varepsilon_j|).$$

According to this and the proof of (7d) we have

$$\hat{r}_j \leq \frac{2\rho_j}{|u_j|^2 - \rho_j^2} = O(r\varepsilon_j^2),$$

which implies the lemma. \square

In the sequel we will derive a qualitative estimate of \hat{E} . Here, only using $()^{I_1}$ or $()^{I_2}$ in (5) will give Nourain's method for the centers which is of fourth order. Hence we will obtain $\hat{E} = O(E^4)$ in this case. To include applications of $()^{-1}$ in (5) we firstly estimate the distance of the midpoints.

Lemma 4.3. *Using the notations from Lemma 4.1 and neglecting the iteration index ν there holds:*

(i) *Let α be equal to 1 if $INV_1 = ()^{-1}$ and 0 otherwise (i.e. $INV_1 = ()^{I_1}$ or $()^{I_2}$). Then, for all $j \in \{1, \dots, m\}$,*

$$|\text{mid}(INV_1(W_j)) - 1/u_j| = \alpha \cdot O(E^3R^2).$$

(ii) *Let β be equal to 1 if $INV_2 = ()^{-1}$ and 0 otherwise. Then, for all $j, k \in \{1, \dots, m\}$, $j \neq k$.*

$$|\text{mid}(INV_2(z_j - Z_k + N_k)) - 1/v_{jk}| = \beta \cdot O(R^2).$$

Proof. For $INV_i = ()^{I_1}$ or $INV_i = ()^{I_2}$ the assertions are trivial ($\alpha = 0$ or $\beta = 0$). Let us assume $INV_i = ()^{-1}$. Then, cf. Section 1,

$$\left| \text{mid}(W_j^{-1}) - \frac{1}{u_j} \right| = \frac{1}{|u_j|} \frac{\rho_j^2/|u_j|^2}{1 - \rho_j^2/|u_j|^2}.$$

Using $|u_j \varepsilon_j| = O(1)$ and $\rho_j = O(R)$ assertion (i) follows.

Similarly,

$$\left| \text{mid}(INV_2(z_j - Z_k + N_k)) - \frac{1}{v_{jk}} \right| \leq \frac{1}{|v_{jk}|} \frac{R^2/|v_{jk}|^2}{1 - R^2/|v_{jk}|^2} = O(R^2),$$

which proves (ii). \square

Together with Lemma 4.2, the following lemma will be the tool for proving Theorem 3.2.

Lemma 4.4. *Using the notations from Lemma 4.1 and Lemma 4.3 and neglecting the iteration index ν there holds*

$$\hat{E} = O(E^4) + \alpha O(E^3R^2) + \beta O(E^2R^2).$$

Proof. Using the above formula for N_j and Lemma 4.3(ii) there holds, for any $j \in \{1, \dots, m\}$,

$$\begin{aligned} \text{mid}(W_j) &= \frac{1}{\varepsilon_j} + \frac{1}{\mu_j} \sum_{k=1, k \neq j}^m \mu_k \left(\frac{1}{z_j - \zeta_k} - \text{mid}(INV_2(z_j - Z_k + N_k)) \right) \\ &= \frac{1}{\varepsilon_j} + \frac{1}{\mu_j} \sum_{k=1, k \neq j}^m \mu_k \left(\frac{-z_k + N_k + \zeta_k}{(z_j - \zeta_k)v_{jk}} + \beta \cdot O(R^2) \right). \end{aligned}$$

Considering the estimate of the error of Newton–Raphson's method in the proof of Lemma 4.1, we see its quadratic convergence

$$|z_j - N_j - \zeta_j| = O(E^2).$$

Thus,

$$u_j = \text{mid}(W_j) = \frac{1}{\varepsilon_j} + O(E^2) + \beta O(R^2).$$

Therefore and by Lemma 4.3(i),

$$\begin{aligned} \hat{\varepsilon}_j &= \varepsilon_j - \text{mid}(\text{INV}_1(W_j)) \\ &= \varepsilon_j + \frac{1}{1/\varepsilon_j + O(E^2) + \beta O(R^2)} + \alpha O(E^3 R^2) \\ &= \varepsilon_j \cdot (1 - (1 + O(E^3) + \beta O(ER^2))^{-1}) + \alpha O(E^3 R^2) \\ &= O(E^4) + \beta O(R^2 E^2) + \alpha O(E^3 R^2), \end{aligned}$$

which proves the lemma. \square

To determine the R-order of convergence of the interval method (5) in the case when $\text{INV}_2 = ()^{-1}$, we need the following result where K_ν is called a convergence factor if the sequence (K_ν) is bounded. The result is known in the theory of iterative processes (and covered e.g. by [5]), so we may omit its elementary proof.

Lemma 4.5. *Let (s_ν) be a sequence of positive numbers tending towards zero such that $s_{\nu+2} \leq K_\nu s_{\nu+1}^p s_\nu^q$, where K_ν is the convergence factor and p and q being natural numbers. Then the R-order of (s_ν) is at least $\frac{1}{2}(p + \sqrt{p^2 + 4q})$.*

After these preliminaries we are now in the position to determine the R-order of convergence of method (5).

Proof of Theorem 3.2. Using the notations of Lemma 4.3 we have to prove that the sequence (R_ν) has R-order $\frac{1}{2}(3 + \sqrt{17}) \approx 3.562$ if $\beta = 1$ and 4 if $\beta = 0$. We may assume that (R_ν) converges towards zero and the conclusions of Lemmas 4.2 and 4.4 hold. In addition, we assume without loss of generality that $E_0 \leq R_0 < 1$ is sufficiently small (otherwise we choose ν_0 sufficiently large and consider the sequence $(R_\nu)_{\nu \geq \nu_0}$).

Firstly, assume $\beta = 0$. Then, according to Lemmas 4.2 and 4.4 there exist constants A_ν, B_ν , and D_ν depending on ν such that for any $\nu = 0, 1, 2, \dots$

$$E_{\nu+1} \leq A_\nu E_\nu^4 + B_\nu E_\nu^3 R_\nu^2, \quad R_{\nu+1} \leq D_\nu R_\nu E_\nu^2.$$

Assume that $D_\nu \geq A_\nu, B_\nu \geq 4$. We note that the constants A_ν, B_ν , and D_ν have the role of convergence factors. Let $e_0 := E_0, r_0 := R_0$, and define for any $\nu = 0, 1, 2, \dots$

$$e_{\nu+1} := A_\nu e_\nu^4 + B_\nu e_\nu^3 r_\nu^2, \tag{8a}$$

$$r_{\nu+1} := D_\nu r_\nu e_\nu^2, \tag{8b}$$

where A_ν, B_ν , and D_ν are positive constants depending on ν such that $D_\nu \geq A_\nu, B_\nu \geq 4$.

By induction, $E_\nu \leq e_\nu$ and $R_\nu \leq r_\nu$. Since R_0 may be sufficiently small, we may assume that the sequences (e_ν) and (r_ν) converge towards zero. It suffices to prove that the Q-order of (r_ν) is at least four. Therefore, let

$$\tau_\nu := \sqrt{\frac{r_{\nu+1}}{D_\nu r_\nu^4}}, \quad \nu = 0, 1, 2, \dots \tag{9}$$

Now it remains to prove that (τ_ν) is bounded. By definition of $r_{\nu+1}$,

$$\tau_\nu = \frac{e_\nu}{r_\nu^{3/2}} \quad \text{i.e.} \quad e_\nu = \tau_\nu r_\nu^{3/2}, \quad \nu = 0, 1, 2, \dots$$

Using $e_\nu = \tau_\nu r_\nu^{3/2}$ and $e_{\nu+1} = \tau_{\nu+1} r_{\nu+1}^{3/2}$ on both sides of equation (8a) we obtain

$$\tau_{\nu+1} r_{\nu+1}^{3/2} = A_\nu \tau_\nu^4 r_\nu^6 + B_\nu \tau_\nu^3 r_\nu^{13/2}. \tag{10}$$

From (9), we get $r_{\nu+1} = D_\nu \tau_\nu^2 r_\nu^4$ which will be substituted on the left-hand side of (10) to obtain after some calculations

$$\tau_{\nu+1} = \frac{A_\nu \tau_\nu + B_\nu r_\nu^{1/2}}{D_\nu^{3/2}} \leq \frac{\tau_\nu + r_\nu^{1/2}}{2}.$$

Since (r_ν) tends towards zero it is bounded by K' , say. Letting $K'' := \max\{\sqrt{K'}, \tau_0\}$ we find for all $\nu = 0, 1, 2, \dots$

$$\tau_{\nu+1} \leq \frac{1}{2}(\tau_\nu + K'').$$

By induction we check that $\tau_\nu \leq K''$ which proves Theorem 3.2 for $\beta = 0$.

Consider now the case $\beta = 1$. As in the first part of the proof, according to Lemmas 4.2 and 4.4 there exist constants A_ν, B_ν, C_ν , and D_ν such that for any $\nu = 0, 1, 2, \dots$

$$E_{\nu+1} \leq A_\nu E_\nu^4 + B_\nu E_\nu^3 R_\nu^2 + C_\nu E_\nu^2 R_\nu^2,$$

$$R_{\nu+1} \leq D_\nu R_\nu E_\nu^2.$$

Assume that $D_\nu \geq A_\nu, B_\nu, C_\nu \geq 4$. Again, let $e_0 := E_0 \leq r_0 := R_0$ and define for any $\nu = 0, 1, 2, \dots$

$$e_{\nu+1} := A_\nu e_\nu^4 + B_\nu e_\nu^3 r_\nu^2 + C_\nu e_\nu^2 r_\nu^2, \tag{11a}$$

$$r_{\nu+1} := D_\nu r_\nu e_\nu^2. \tag{11b}$$

By induction we prove that $E_\nu \leq e_\nu$ and $R_\nu \leq r_\nu$ so it remains to estimate the R-order of (r_ν) . Again, we may assume that R_0 is sufficiently small and (e_ν) and (r_ν) converge towards zero. Moreover, we may assume that $e_\nu \leq \frac{1}{3}$ and $r_\nu \leq \frac{1}{3}$. Using the above definitions we obtain

$$\frac{e_{\nu+1}}{r_{\nu+1}} = \frac{A_\nu e_\nu^4 + B_\nu e_\nu^3 r_\nu^2 + C_\nu e_\nu^2 r_\nu^2}{D_\nu r_\nu e_\nu^2} \leq \frac{e_\nu^2 + e_\nu r_\nu^2 + r_\nu^2}{r_\nu} \leq \frac{4}{9} \cdot \frac{e_\nu}{r_\nu} + \frac{1}{3}.$$

From this, we conclude by induction that $e_\nu \leq r_\nu$ for any $\nu = 0, 1, 2, \dots$, which implies (by definition of $r_{\nu+1}$ (11b)) that

$$r_{\nu+1} \leq D_\nu r_\nu^3, \quad \nu = 0, 1, 2, \dots \tag{12}$$

Taking (11b) to obtain

$$e_\nu = \sqrt{\frac{r_{\nu+1}}{D_\nu r_\nu}}$$

and substituting this (and the related expression for $e_{\nu+1}$) in the defining relation for $e_{\nu+1}$ (11a) we get

$$\left(\frac{r_{\nu+2}}{D_{\nu+1} r_{\nu+1}}\right)^{1/2} = \frac{A_\nu r_{\nu+1}^2}{D_\nu^2 r_\nu^2} + B_\nu \left(\frac{r_{\nu+1}}{D_\nu r_\nu}\right)^{3/2} r_\nu^2 + C_\nu \frac{r_{\nu+1} r_\nu}{D_\nu}.$$

Hence, using (12) and the inequalities $D_\nu \geq A_\nu, B_\nu, C_\nu \geq 4$, we find

$$r_{\nu+2} \leq D_{\nu+1} r_{\nu+1}^3 r_\nu^2 \left[\frac{1}{D_\nu} \cdot \frac{r_{\nu+1}}{r_\nu^3} + \frac{1}{D_\nu^{1/2}} \left(\frac{r_{\nu+1}}{r_\nu^3}\right)^{1/2} r_\nu + \frac{C_\nu}{D_\nu} \right]^2 \leq \frac{49}{9} D_{\nu+1} r_{\nu+1}^3 r_\nu^2.$$

Then, Lemma 4.5 (for $p = 3$ and $q = 2$) proves that (r_ν) has R -order $\frac{1}{2}(3 + \sqrt{17})$. \square

Remark. Using the conclusions of Lemmas 4.2 and 4.4 the theorem can also be proved using [5].

5. Numerical results

From the convergence analysis of the new algorithms presented in Sections 3 and 4, we see that these algorithms have a very fast convergence. The values of their convergence order 3.562 and 4 should be regarded as asymptotical ones meaning that the notified speed of convergence can be realized after several iterative steps. Numerical examples have shown that such situation most frequently begins with the third or, at best, the second iterative step. Thus, the presented methods with Schröder's correction are the most powerful when at least three iterations are applied. This fact is, as in the case of all iterative methods with very fast convergence, a slight disadvantage in a certain degree because floating-point arithmetic of very high accuracy should be employed (usually, quadruple precision arithmetic).

The iterative formula (5) enables the construction of several methods depending on the choice of the type of inversion. We have considered three methods choosing INV_1 in (5) to be the exact inversion (i.e. $\text{INV}_1(W) = W^{-1}$) and taking $\text{INV}_2 \in \{()^{-1}, ()^{I_1}, ()^{I_2}\}$. These methods are referred to as M_I , M_{I_1} , and M_{I_2} and displayed below by the following formulas:

$$\hat{Z}_j = z_j - \left(\frac{1}{N_j} - \frac{1}{\mu_j} \sum_{k=1, k \neq j}^m \mu_k (z_j - Z_k + N_k)^{-1} \right)^{-1}, \quad (M_I)$$

$$\hat{Z}_j = z_j - \left(\frac{1}{N_j} - \frac{1}{\mu_j} \sum_{k=1, k \neq j}^m \mu_k (z_j - Z_k + N_k)^{I_1} \right)^{-1}, \quad (M_{I_1})$$

$$\hat{Z}_j = z_j - \left(\frac{1}{N_j} - \frac{1}{\mu_j} \sum_{k=1, k \neq j}^m \mu_k (z_j - Z_k + N_k)^{I_2} \right)^{-1}, \quad (M_{I_2})$$

for $j = 1, \dots, m$. From the definition of the inverse $()^{I_1}$ we observe that the evaluation of $|z|$ is required. But this is very costly because the evaluation of $|z| = |x + iy| = \sqrt{x^2 + y^2}$ needs too much extra time. Consequently, the method M_{I_1} has a poor computational efficiency (see Section 6). For this reason we are forced to introduce the inversion $()^{I_2}$. Namely, since $|z| > r$ we have

$$\frac{r}{|z|(|z| - r)} = \frac{r(1 + r/|z|)}{|z|^2 - r^2} < \frac{2r}{|z|^2 - r^2}.$$

Thus, $\text{rad}\{z, r\}^{I_2} = 2 \text{rad}\{z, r\}^{-1}$ which is negligibly increased (compared to the exact inversion) if r is very small (for example, after the first or the second iteration). However, the tested numerical examples have shown that this increase of the radius of the inverse disc causes certain difficulties if the initial discs are not sufficiently small. This problem can be solved partly using the approximation $x < \frac{1}{2}(1 + x^2)$ ($0 < x = r/|z| < 1$) to obtain

$$\frac{r}{|z|(|z| - r)} < \frac{r\left(1 + \frac{1}{2}(1 + r^2/|z|^2)\right)}{|z|^2 - r^2} = \frac{r\left(\frac{3}{2} + \frac{1}{2}r^2/|z|^2\right)}{|z|^2 - r^2} < \frac{2r}{|z|^2 - r^2}.$$

The new type of inversion will be denoted by

$$\{z, r\}^{\hat{I}_2} = \left\{ \frac{1}{z}, \frac{r\left(\frac{3}{2} + \frac{1}{2}r^2/|z|^2\right)}{|z|^2 - r^2} \right\}.$$

Since r^2 and $|z|^2$ are known, new calculations are not necessary. We will denote the corresponding method with $M_{\hat{I}_2}$.

Although the above decrease of the radius of the inverse disc is small, it is often sufficient to provide very good results. Moreover, it is sufficient to apply the inverse $()^{\hat{I}_2}$ only in the first iterative step. Such a combined method will be referred to as $M_{\hat{I}_2, I_2}$.

The proposed three methods, as well as the third-order interval method (2) for the comparison purpose, have been tested in solving many polynomial equations. The programs have been written in FORTRAN language and implemented on the computer VAX 3400 in quadruple precision arithmetic (about 33 significant decimal digits). For demonstration, we present the following example.

Example. We have used the presented algorithms for finding the zeros of a polynomial

$$\begin{aligned} P(z) = & z^9 + (-2 + 3i)z^8 + (48 - 6i)z^7 + (-94 + 152i)z^6 + (522 - 298i)z^5 \\ & + (-950 + 1974i)z^4 + (-1400 - 3650i)z^3 + (3750 + 1200i)z^2 \\ & + (-1875 + 1250i)z - 625i. \end{aligned}$$

The exact zeros of P are $\zeta_1 = 1$, $\zeta_2 = -i$, $\zeta_3 = -5i$, and $\zeta_4 = 5i$ with the multiplicities $\mu_1 = 2$, $\mu_2 = 3$, $\mu_3 = 2$, and $\mu_4 = 2$. As the initial inclusion discs we have taken the circular regions

$$\begin{aligned} Z_1^{(0)} &= \{1.1 + 0.2i, 0.9\}, & Z_2^{(0)} &= \{0.2 - 0.8i, 0.9\}, \\ Z_3^{(0)} &= \{-0.6 - 4.4i, 0.9\}, & Z_4^{(0)} &= \{-0.6 + 4.4i, 0.9\}. \end{aligned}$$

The maximal radii $r^{(\nu)} = \max_{j=1, \dots, m} r_j^{(\nu)}$ ($\nu = 1, 2, 3$) for all the methods are shown in Table 1.

Table 1
Maximal radii of inclusion discs

	Method (2)	M_I	M_{I_1}	M_{I_2}	$M_{\hat{I}_2}$	$M_{\hat{I}_2, I_2}$
$r^{(1)}$	1.16(-1)	1.25(-1)	2.44(-1)	3.33(-1)	2.35(-1)	2.35(-1)
$r^{(2)}$	9.55(-4)	3.78(-5)	5.19(-4)	3.54(-3)	7.47(-4)	9.96(-4)
$r^{(4)}$	4.35(-13)	3.61(-17)	5.18(-16)	1.24(-12)	1.5(-15)	3.51(-15)

From the last two columns we observe that the application of the new inversion $()^{\hat{I}_2}$ produces considerably improved discs compared to M_{I_2} . The interval method M_{I_2} shows a fast convergence after the second iteration but the discs $Z_j^{(2)}$ are not sufficiently contracted. Finally, the interval method M_I which uses only the exact inversions generates good results. The above conclusions are also valid for most of the tested polynomials.

Remarks.

- (i) In order to realize self-verifying methods additional calculations are required which enlarge the discs as well. This could be a disadvantage compared with rectangular complex interval arithmetic.
- (ii) As seen in the proof of Theorem 3.2, the increased convergence of our methods is caused by using “better midpoints”. Hence, the perturbation of the midpoints (e.g. due to roundoff errors) results in a less increased convergence speed.

6. Computational efficiency

An estimate of the computational efficiency of the proposed interval methods can be sensibly carried out using the *coefficient of efficiency* $E_n(M)$ given by

$$E_n(M) = r(n)^{c/\theta(n)}$$

(see [14, Chapter 6]). Here $r(n)$ is the R-order of convergence of the method (M), $\theta(n)$ is the computational cost per iteration, and c is a normalization constant. $\theta(n)$ is proportional to the total number of the basic arithmetic operations taken with certain *operation weight* depending on CPU time.

We have calculated the coefficient of efficiency for four computing machines and for $n = 3, 4, 5, \dots, 15$ considering the case of simple zeros. For comparison purpose we have included Gargantini's method (2) referred to as M_G which used to be the most efficient method for the simultaneous inclusion of polynomial zeros (see [14, Chapter 6]).

First, we have used that the order of convergence of interval methods M_{I_1} , M_{I_2} , $M_{\hat{I}_2}$ and $M_{\hat{I}_2, I_2}$ is 4, and equal to 3.562 for the interval method M_I . The rating of the considered methods was the same for all four computers and all $n = 3(1)15$. We have obtained the following rank:

1.	2.	3.	4.	5.	6.
M_{I_2}	$M_{\hat{I}_2, I_2}$	M_I	$M_{\hat{I}_2}$	M_G	M_{I_1}

We emphasize that the above ordering is of a theoretical importance because we used the asymptotical values of the order of convergence. As discussed in the previous section, the convergence speed is somewhat smaller than the theoretical one. Calculating with average values of convergence orders determined by considering the first three iterations and taking $r(n) = 3$ for the method M_G , we have obtained the following rating:

1.–2.	3.	4.	5.	6.
M_{I_2}	$M_{\hat{I}_2, I_2}$	M_G	$M_{\hat{I}_2}$	M_{I_1}
M_I				

The values of the coefficient of efficiency of the interval methods M_{I_2} and M_I are very close so that both methods are ranged equally.

According to the above lists we can conclude that the interval methods M_{I_2} , M_I and $M_{\hat{I}_2, I_2}$ are more efficient than M_G . Consequently, these methods become the most efficient methods for the simultaneous inclusion of polynomial zeros.

7. References

- [1] O. Aberth, Iteration methods for finding all zeros of a polynomial simultaneously, *Math. Comp.* 27 (1973) 339–344.
- [2] ACRITH-XCS IBM, High accuracy arithmetic-extended scientific computation, ACRITH-XSC Language Reference, SC33-6462-00, IBM Corporation (1990).
- [3] G. Alefeld and J. Herzberger, *Introduction to Interval Computation* (Academic Press, New York, 1983).
- [4] W. Börsch-Supan, A posteriori error bounds for the zeros of polynomials, *Numer. Math.* 5 (1963) 380–398.
- [5] W. Burmeister and J.W. Schmidt, On the R-order of coupled sequences arising in single-step type methods, *Numer. Math.* 53 (1988) 653–661.
- [6] L.W. Ehrlich, A modified Newton method for polynomials, *Comm. ACM* 10 (1967) 107–108.
- [7] I. Gargantini, Further applications of circular arithmetic: Schroeder-like algorithms with error bounds for finding zeros of polynomials, *SIAM J. Numer. Math.* 15 (1978) 497–510.
- [8] I. Gargantini and P. Henrici, Circular arithmetic and the determination of polynomial zeros, *Numer. Math.* 18 (1972) 305–320.
- [9] P. Henrici, *Applied and Computational Complex Analysis, Vol. I* (Wiley, New York, 1974).
- [10] R. Klatté, U. Kulisch, M. Neaga, D. Ratz and C. Ullrich, *PASCAL-XSC, Sprachbeschreibung mit Beispielen* (Springer, Berlin, 1991).
- [11] V.H. Maehly, Zur iterativen Auflösung algebraischer Gleichungen, *Z. Angew. Math. Phys.* 5 (1954) 260–263.
- [12] A.W.M. Nourain, An improvement on two iteration methods for simultaneous determination of the zeros of a polynomial, *Internat. J. Comput. Math.* 6 (1977) 241–252.
- [13] J.M. Ortega and W.C. Rheinboldt, *Iterative Solution of Nonlinear Equations in Several Variables* (Academic Press, New York, 1970).
- [14] M.S. Petković, *Iterative Methods for Simultaneous Inclusion of Polynomial Zeros* (Springer, Berlin, 1989).