

## ON LOCATING CLUSTERS OF ZEROS OF ANALYTIC FUNCTIONS \*

P. KRAVANJA<sup>1</sup>, T. SAKURAI<sup>2</sup>, and M. VAN BAREL<sup>1</sup> †

<sup>1</sup>*Department of Computer Science, Katholieke Universiteit Leuven  
Celestijnenlaan 200 A, B-3001 Heverlee, Belgium.  
email: Peter.Kravanja@na-net.ornl.gov, Marc.VanBarel@cs.kuleuven.ac.be*

<sup>2</sup>*Institute of Information Sciences and Electronics, University of Tsukuba  
Tsukuba 305, Japan. email: sakurai@is.tsukuba.ac.jp*

### Abstract.

Given an analytic function  $f$  and a Jordan curve  $\gamma$  that does not pass through any zero of  $f$ , we consider the problem of computing *all* the zeros of  $f$  that lie inside  $\gamma$ , together with their respective multiplicities. Our principal means of obtaining information about the location of these zeros is a certain symmetric bilinear form that can be evaluated via numerical integration along  $\gamma$ . If  $f$  has one or several clusters of zeros, then the mapping from the ordinary moments associated with this form to the zeros and their respective multiplicities is very ill-conditioned. We present numerical methods to calculate the centre of a cluster and its weight, i.e., the arithmetic mean of the zeros that form a certain cluster and the total number of zeros in this cluster, respectively. Our approach relies on formal orthogonal polynomials and rational interpolation at roots of unity. Numerical examples illustrate the effectiveness of our techniques.

*AMS subject classification:* Primary 65H05; Secondary 65E05.

*Key words:* Zeros of analytic functions, clusters of zeros, logarithmic residue integrals, formal orthogonal polynomials, rational interpolation.

### 1 Introduction.

Let  $W$  be a simply connected region in  $\mathbb{C}$ ,  $f : W \rightarrow \mathbb{C}$  analytic in  $W$ , and  $\gamma$  a positively oriented Jordan curve in  $W$  that does not pass through any zero of  $f$ . We consider the problem of computing *all* the zeros of  $f$  that lie in the interior of  $\gamma$ , together with their respective multiplicities. Our principal means of obtaining information about the location of these zeros is a certain symmetric bilinear form that can be evaluated via numerical integration along  $\gamma$ . If  $f$  has one or several clusters of zeros, then the mapping from the ordinary moments

---

\*Received July 1993. Revised January 1994.

†The first author was supported by a grant from the Flemish Institute for the Promotion of Scientific and Technological Research in the Industry (IWT). This work is part of the projects #G.0261.96 “Counting and Computing all Isolated Solutions of Systems of Nonlinear Equations” and #G.0278.97 “Orthogonal Systems and their Applications” funded by the Fund for Scientific Research, Flanders.

associated with this form to the zeros and their respective multiplicities is very ill-conditioned. We will present numerical methods to calculate the centre of a cluster and its weight, i.e., the arithmetic mean of the zeros that form a certain cluster and the total number of zeros in this cluster, respectively. This information enables one to zoom into a certain cluster: its zeros can be calculated separately from the other zeros of  $f$ . By shifting the origin in the complex plane to the centre of a certain cluster, its zeros become better relatively separated, which is appropriate in floating-point arithmetic and reduces the ill-conditioning.

Our approach to the problem of computing zeros of analytic functions can be seen as a continuation of the pioneering work by Delves and Lyness [17]. Let  $N$  denote the total number of zeros of  $f$  that lie in the interior of  $\gamma$ , i.e., the number of zeros where each zero is counted according to its multiplicity. Suppose that  $N > 0$ . Let  $n$  be the number of mutually distinct zeros of  $f$  that lie in the interior of  $\gamma$ . Let  $z_1, \dots, z_n$  be these zeros and  $\nu_1, \dots, \nu_n$  their respective multiplicities. An easy calculation shows that  $z_k$  is a simple pole of  $f'/f$  with residue  $\nu_k$  for  $k = 1, \dots, n$ . It follows that

$$(1.1) \quad N = \frac{1}{2\pi i} \int_{\gamma} \frac{f'(z)}{f(z)} dz$$

and thus  $N$  can be calculated via numerical integration. Methods for the determination of zeros of analytic functions that are based on the numerical evaluation of integrals are called *quadrature methods*. A review of such methods was given by Ioakimidis [29]. Delves and Lyness considered the integrals

$$s_p := \frac{1}{2\pi i} \int_{\gamma} z^p \frac{f'(z)}{f(z)} dz, \quad p = 0, 1, 2, \dots$$

The residue theorem implies that the  $s_p$ 's are equal to the *Newton sums* of the unknown zeros,

$$(1.2) \quad s_p = \sum_{k=1}^n \nu_k z_k^p, \quad p = 0, 1, 2, \dots$$

In what follows we will assume that all the  $s_p$ 's that are needed have been calculated. In particular, we will assume that the value of  $N = s_0$  is known.

Delves and Lyness considered the following monic polynomial of degree  $N$ :

$$P_N(z) := \prod_{k=1}^n (z - z_k)^{\nu_k} =: z^N + \sigma_1 z^{N-1} + \dots + \sigma_N.$$

They called  $P_N(z)$  the *associated polynomial* for the interior of  $\gamma$ . Its coefficients can be calculated via Newton's identities.

THEOREM 1.1 (NEWTON'S IDENTITIES).

$$\begin{aligned} s_1 + \sigma_1 &= 0, \\ s_2 + s_1 \sigma_1 + 2 \sigma_2 &= 0, \\ &\vdots \\ s_N + s_{N-1} \sigma_1 + \dots + s_1 \sigma_{N-1} + N \sigma_N &= 0. \end{aligned}$$

PROOF. An elegant proof was given by Carpentier and Dos Santos [16].  $\square$

In this way they reduced the problem to the easier problem of computing the zeros of a polynomial. Unfortunately, the map from the Newton sums  $s_1, \dots, s_N$  to the coefficients  $\sigma_1, \dots, \sigma_N$  is usually ill-conditioned. Also, the polynomials that arise in practice may be such that small changes in the coefficients produce much larger changes in some of the zeros. This ill-conditioning of the map between the coefficients of a polynomial and its zeros was investigated by Wilkinson [52]. The location of the zeros determines their sensitivity to perturbations of the coefficients. Multiple zeros and very close zeros are extremely sensitive, but even a succession of moderately close zeros can result in severe ill-conditioning. Wilkinson states that ill-conditioning in polynomials cannot be overcome without, at some stage of the computation, resorting to high precision arithmetic.

If  $f$  has many zeros in the interior of  $\gamma$ , then the associated polynomial is of high degree and could be very ill-conditioned. Therefore, if  $N$  is large, one has to calculate the coefficients  $\sigma_1, \dots, \sigma_N$ , and thus the integrals  $s_1, \dots, s_N$ , very accurately. To avoid the use of high precision arithmetic and to reduce the number of integrand evaluations needed to approximate the  $s_p$ 's, Delves and Lyness suggested to construct and solve the associated polynomial only if its degree is smaller than or equal to a preassigned number  $M$ . Otherwise, the interior of  $\gamma$  is subdivided or covered with a finite covering and the smaller regions are treated in turn. The choice of  $M$  involves a trade-off. If  $M$  is increased, then fewer regions have to be scanned. However, if  $M$  is chosen too large, then the resulting associated polynomial may be ill-conditioned. Delves and Lyness chose  $M = 5$ .

Botten, Craig and McPhedran [12] made a Fortran 77 implementation of the method of Delves and Lyness.

Instead of using Newton's identities to construct the associated polynomial, Li [39] considered (1.2) as a system of polynomial equations. He used a homotopy continuation method to solve this system.

What is wrong with these approaches, in our opinion, is that they consider the wrong set of unknowns. One should consider the mutually distinct zeros and their respective multiplicities *separately*. This is the approach that we will follow.

This paper is organized as follows. In Section 2 we introduce a symmetric bilinear form  $\langle \cdot, \cdot \rangle$  that is related to  $f$  and  $\gamma$  and we study the polynomials that are mutually orthogonal with respect to this form. These polynomials are called *formal orthogonal polynomials* (FOPs). The problem of computing  $z_1, \dots, z_n$  is shown to be equivalent to that of computing the zeros of the  $n$ th degree FOP. We will show how zeros of FOPs can be calculated by solving generalized eigenvalue problems. The value of  $n$  will be determined indirectly. Once  $n$  and  $z_1, \dots, z_n$  have been found, the problem becomes linear and the multiplicities  $\nu_1, \dots, \nu_n$  can be computed by solving a Vandermonde system. In Section 3 we suppose that  $z_1, \dots, z_n$  can be grouped into  $m$  clusters. We introduce a symmetric bilinear

form  $\langle \cdot, \cdot \rangle_m$  that is related to the centres and the weights of the clusters, and we show that the form  $\langle \cdot, \cdot \rangle_m$  approximates  $\langle \cdot, \cdot \rangle$ . In Section 4 we present an algorithm to compute FOPs in their product representation, i.e., the polynomials are represented in terms of their zeros. In Section 6 we attack our problem in an entirely different way, based on rational interpolation at roots of unity. However, we will show that the denominator polynomials of the interpolants are FOPs with respect to a certain form  $\langle \langle \cdot, \cdot \rangle \rangle$  and that, under certain conditions, the forms  $\langle \cdot, \cdot \rangle$  and  $\langle \langle \cdot, \cdot \rangle \rangle$  coincide. This enables us to calculate the FOPs with respect to  $\langle \cdot, \cdot \rangle$  in a different way. Numerical examples are presented in Sections 5 and 7.

A Fortran 90 implementation of the approach for computing zeros of analytic functions presented in this paper will appear in [38]. This package cannot be used for locating clusters, only for computing zeros and their respective multiplicities.

REMARK 1.1. Our aim is to present techniques for computing zeros of analytic functions that give very accurate results. As the reader will notice, if we have a choice between several options for a certain part of an algorithm, then we will always choose the option that, in our experience, gives the most accurate results, even if it is the most expensive (though still within the limits of what is reasonable, of course) in terms of number of floating-point operations. The emphasis lies on accuracy.

REMARK 1.2. If  $f$  has one or more zeros on the curve  $\gamma$ , then  $f'/f$  has a pole on  $\gamma$  and the integral (1.1) is no longer well defined. This is the reason why we do not allow  $f$  to have zeros on  $\gamma$ .

Let us briefly mention a number of related approaches.

Petković *et al.* [44, 45, 46] presented simultaneous iterative methods for computing zeros of analytic functions. These algorithms can be used only if the zeros of  $f$  are known to be simple and if sufficiently accurate initial approximations are available. Atanassova [9] considered the case of multiple zeros but assumed that the multiplicities are known in advance.

In a number of papers and short notes, Anastasselou and Ioakimidis [1, 2, 3, 4, 5, 30, 31, 32, 33] considered the problem of computing zeros of sectionally analytic functions (i.e., functions that are analytic except for a finite number of discontinuity arcs). They proposed variations and generalizations of the method of Burniston and Siewert [14], which is based on the theory of Riemann–Hilbert boundary value problems (cf. Gakhov [20]). The authors focused on the function  $\alpha + \beta z - z \tanh^{-1}(1/z)$  where  $\alpha, \beta \in \mathbb{C}$  are parameters. This function appears in the theory of ferromagnetism. It has the discontinuity interval  $[-1, 1]$ . Already for this example, the approach of Anastasselou and Ioakimidis requires a lot of specific analytical calculations. Therefore, it is unclear how their method could lead to a ‘black box’ algorithm that can handle arbitrary functions. Also, although multiple zeros are not a problem, their approach cannot calculate multiplicities.

Yakoubsohn [53] proposed an exclusion method for computing zeros of analytic functions. Unfortunately, his exclusion function is difficult to evaluate and multiple zeros require special treatment. He applied his algorithm only to polynomials. See also Ying and Katz [54].

Specifically for clusters of polynomial zeros, let us mention that Hribernig and Stetter [28] worked on detection and validation of clusters of zeros whereas Kirrinnis [35] studied Newton's iteration towards a cluster. See also Neumaier [42].

In [36] Kravanja, Cools and Haegemans used a multidimensional version of the integral formula (1.2) to solve systems of analytic equations. In [37] Kravanja, Van Barel and Haegemans used an approach based on formal orthogonal polynomials to compute zeros and poles of meromorphic functions. See also [49].

## 2 Formal orthogonal polynomials.

Let  $\mathcal{P}$  be the linear space of polynomials with complex coefficients. We define a symmetric bilinear form

$$\langle \cdot, \cdot \rangle : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{C}$$

by setting

$$(2.1) \quad \langle \phi, \psi \rangle := \frac{1}{2\pi i} \int_{\gamma} \phi(z) \psi(z) \frac{f'(z)}{f(z)} dz = \sum_{k=1}^n \nu_k \phi(z_k) \psi(z_k)$$

for any two polynomials  $\phi, \psi \in \mathcal{P}$ . Note that  $\langle \cdot, \cdot \rangle$  can be evaluated via numerical integration. Let  $s_p := \langle 1, z^p \rangle$  for  $p = 0, 1, 2, \dots$ . These ordinary moments are equal to the Newton sums of the unknown zeros,

$$s_p = \sum_{k=1}^n \nu_k z_k^p, \quad p = 0, 1, 2, \dots$$

In particular,  $s_0 = N$ . Let  $H_k$  be the  $k \times k$  Hankel matrix

$$H_k := \begin{bmatrix} s_0 & s_1 & \cdots & s_{k-1} \\ s_1 & & \ddots & \vdots \\ \vdots & \ddots & & \vdots \\ s_{k-1} & \cdots & \cdots & s_{2k-2} \end{bmatrix}$$

for  $k = 1, 2, \dots$ . Note that the form  $\langle \cdot, \cdot \rangle$  is completely determined by the sequence of moments  $(s_p)_{p \geq 0}$ . In Section 1 we mentioned that the problem of computing all the zeros of  $f$  that lie inside  $\gamma$  is ill-conditioned in case  $f$  has one or several clusters of zeros. Mathematically, this refers to the conditioning of the mapping  $(s_p)_{p \geq 0} \mapsto (z_1, \dots, z_n, \nu_1, \dots, \nu_n)$ .

A monic polynomial  $\varphi_t$  of degree  $t \geq 0$  that satisfies

$$(2.2) \quad \langle z^k, \varphi_t(z) \rangle = 0, \quad k = 0, 1, \dots, t-1,$$

is called a *formal orthogonal polynomial* (FOP) of degree  $t$ . (Observe that condition (2.2) is void for  $t = 0$ .) The adjective *formal* emphasizes the fact that, in general, the form  $\langle \cdot, \cdot \rangle$  does not define a true inner product. An important

consequence of this fact is that, in contrast to polynomials that are orthogonal with respect to a true inner product, FOPs  $\varphi_t$  need not exist or need not be unique for every degree  $t$ . (For details, see for example [23] and the references cited therein.) If (2.2) is satisfied and  $\varphi_t$  is unique, then  $\varphi_t$  is called a *regular* FOP and  $t$  a *regular index*. If we set

$$\varphi_t(z) =: u_{0,t} + u_{1,t}z + \cdots + u_{t-1,t}z^{t-1} + z^t$$

then condition (2.2) translates into the Yule–Walker system

$$(2.3) \quad \begin{bmatrix} s_0 & s_1 & \cdots & s_{t-1} \\ s_1 & & \ddots & \vdots \\ \vdots & \ddots & & \vdots \\ s_{t-1} & \cdots & \cdots & s_{2t-2} \end{bmatrix} \begin{bmatrix} u_{0,t} \\ u_{1,t} \\ \vdots \\ u_{t-1,t} \end{bmatrix} = - \begin{bmatrix} s_t \\ s_{t+1} \\ \vdots \\ s_{2t-1} \end{bmatrix}.$$

Hence, the regular FOP of degree  $t \geq 1$  exists if and only if the matrix  $H_t$  is nonsingular.

The following theorem characterizes  $n$ , the number of mutually distinct zeros.

**THEOREM 2.1.**  $n = \text{rank } H_{n+p}$  for every nonnegative integer  $p$ . In particular,  $n = \text{rank } H_N$ .

**PROOF.** Let  $p$  be a nonnegative integer. The matrix  $H_{n+p}$  can be written as

$$\begin{aligned} H_{n+p} &= \sum_{k=1}^n \nu_k \begin{bmatrix} 1 & \cdots & z_k^{n+p-1} \\ \vdots & & \vdots \\ z_k^{n+p-1} & \cdots & z_k^{2(n+p)-2} \end{bmatrix} \\ &= \sum_{k=1}^n \nu_k \begin{bmatrix} 1 \\ \vdots \\ z_k^{n+p-1} \end{bmatrix} \begin{bmatrix} 1 & \cdots & z_k^{n+p-1} \end{bmatrix}. \end{aligned}$$

This implies that  $\text{rank } H_{n+p} \leq n$ . However,  $H_n$  is nonsingular. Indeed, one can easily verify that  $H_n$  can be factorized as  $H_n = V_n D_n V_n^T$  where  $V_n$  is the Vandermonde matrix  $V_n := [z_s^r]_{r=0, s=1}^{n-1, n}$  and  $D_n$  is the diagonal matrix  $D_n := \text{diag}(\nu_1, \dots, \nu_n)$ . Therefore  $\text{rank } H_{n+p} \geq n$ . It follows that  $\text{rank } H_{n+p} = n$ .  $\square$

Thus  $H_n$  is nonsingular whereas  $H_t$  is singular for  $t > n$ . Note that  $H_1 = [s_0]$  is nonsingular by assumption. The regular FOP of degree 1 exists and is given by  $\varphi_1(z) = z - \mu$  where

$$\mu := \frac{s_1}{s_0} = \frac{\sum_{k=1}^n \nu_k z_k}{\sum_{k=1}^n \nu_k}$$

is the arithmetic mean of the zeros. Theorem 2.1 implies that the regular FOP  $\varphi_n$  of degree  $n$  exists and tells us also that regular FOPs of degree larger than  $n$  do not exist. The polynomial  $\varphi_n$  is easily seen to be

$$(2.4) \quad \varphi_n(z) = (z - z_1) \cdots (z - z_n).$$

It is the monic polynomial of degree  $n$  that has  $z_1, \dots, z_n$  as simple zeros.

REMARK 2.1. Kronecker's Theorem [43, p. 37] tells us that the infinite Hankel matrix  $H := [s_{p+q}]_{p,q \geq 0}$  has finite rank if and only if its *symbol*, which is defined as the formal Laurent series

$$\frac{s_0}{z} + \frac{s_1}{z^2} + \frac{s_2}{z^3} + \dots,$$

represents a rational function of  $z$ . This is indeed the case. It is easily seen that

$$(2.5) \quad \sum_{k=1}^n \frac{\nu_k}{z - z_k} = \frac{s_0}{z} + \frac{s_1}{z^2} + \frac{s_2}{z^3} + \dots \quad \text{near } z = \infty.$$

In systems theory [34] the problem of reconstructing the rational function that appears in the left-hand side of (2.5) from the sequence of moments  $(s_p)_{p \geq 0}$  is called a *minimal realization problem*. In that context the moments are called *Markov parameters*. Note that the left-hand side of (2.5) is a rational function of type  $[n - 1/n]$  and that its denominator polynomial is given by  $\varphi_n(z)$ .

If  $H_n$  is strongly nonsingular, i.e., if all its leading principal submatrices are nonsingular, then we have a full set  $\{\varphi_0, \varphi_1, \dots, \varphi_n\}$  of regular FOPs.

What happens if  $H_n$  is not strongly nonsingular? By filling up the gaps in the sequence of existing regular FOPs it is possible to define a sequence  $\{\varphi_t\}_{t=0}^\infty$ , with  $\varphi_t$  a monic polynomial of degree  $t$ , such that if these polynomials are grouped into blocks according to the sequence of regular indices, then polynomials belonging to different blocks are orthogonal with respect to  $\langle \cdot, \cdot \rangle$ . More precisely, define  $\{\varphi_t\}_{t=0}^\infty$  as follows. If  $t$  is a regular index, then let  $\varphi_t$  be the regular FOP of degree  $t$ . Else define  $\varphi_t$  as  $\varphi_r \psi_{t,r}$  where  $r$  is the largest regular index less than  $t$  and  $\psi_{t,r}$  is an arbitrary monic polynomial of degree  $t - r$ . In the latter case  $\varphi_t$  is called an *inner polynomial*. If  $\psi_{t,r}(z) = z^{t-r}$  then we say that  $\varphi_t$  is defined *by using the standard monomial basis*. These polynomials  $\{\varphi_t\}_{t=0}^\infty$  can be grouped into blocks. Each block starts with a regular FOP and the remaining polynomials are inner polynomials. Note that the last block has infinite length. The block orthogonality property is expressed by the fact that the Gram matrix  $G_n := [\langle \varphi_r, \varphi_s \rangle]_{r,s=0}^{n-1}$  is block diagonal. The diagonal blocks are nonsingular, symmetric and zero above the main antidiagonal. If all the inner polynomials in a certain block are defined by using the standard monomial basis, then the corresponding diagonal block has Hankel structure. (See Bultheel and Van Barel [13] for more details.)

Theorem 2.1, (2.1) and (2.4) immediately imply the following.

THEOREM 2.2. *Let  $t \geq n$ . Then  $\varphi_t(z_k) = 0$  for  $k = 1, \dots, n$  and  $\langle z^p, \varphi_t(z) \rangle = 0$  for all  $p \geq 0$ .*

Thus the mutually distinct zeros  $z_1, \dots, z_n$  of  $f$  that lie inside  $\gamma$  are among the zeros of the FOP  $\varphi_t(z)$  for all  $t \geq n$ . In Section 4 we will return to the question of how to obtain the value of  $n$  and how to compute FOPs. Theorem 2.1 suggests to determine  $n$  from the singular value decomposition of  $H_N$ . Indeed, theoretically the  $N - n$  smallest singular values of  $H_N$  are equal to zero. In

practice however, this will not be the case, and it may be difficult to determine the rank of  $H_N$  and hence the value of  $n$  in case the gap between the computed approximations for the zero singular values and the nonzero singular values is too small. Instead, the value of  $n$  will be determined indirectly. We will show how FOPs can be computed in their product representation. Therefore the polynomial  $\varphi_n(z)$  immediately leads to the zeros  $z_1, \dots, z_n$ .

Once  $n$  and  $z_1, \dots, z_n$  have been found, the multiplicities  $\nu_1, \dots, \nu_n$  can be calculated by solving the Vandermonde system

$$\begin{bmatrix} 1 & \cdots & 1 \\ z_1 & \cdots & z_n \\ \vdots & & \vdots \\ z_1^{n-1} & \cdots & z_n^{n-1} \end{bmatrix} \begin{bmatrix} \nu_1 \\ \nu_2 \\ \vdots \\ \nu_n \end{bmatrix} = \begin{bmatrix} s_0 \\ s_1 \\ \vdots \\ s_{n-1} \end{bmatrix}.$$

This can be done via the algorithm of Gohberg and Koltracht [22]. This algorithm takes full account of the structure of a Vandermonde matrix and is not only faster but also more accurate than general purpose algorithms such as Gaussian elimination with partial pivoting. It has arithmetic complexity  $\mathcal{O}(n^2)$ .

### 3 Clusters of zeros.

Suppose that the zeros of  $f$  that lie inside  $\gamma$  can be grouped into  $m$  clusters. Let  $I_1, \dots, I_m$  be index sets that define these clusters, and let

$$\mu_j := \sum_{k \in I_j} \nu_k \quad \text{and} \quad c_j := \frac{1}{\mu_j} \sum_{k \in I_j} \nu_k z_k$$

for  $j = 1, \dots, m$ . In other words,  $\mu_j$  is equal to the total number of zeros that form cluster  $j$  (its “weight”) whereas  $c_j$  is equal to the arithmetic mean of the zeros in cluster  $j$  (its “centre of gravity”). We assume that the centres  $c_1, \dots, c_m$  are mutually distinct. For  $k = 1, \dots, n$  we also define  $\zeta_k := z_k - c_j$  if  $k \in I_j$ . From the definition of  $\mu_j$  and  $c_j$  it follows that

$$\sum_{k \in I_j} \nu_k \zeta_k = 0, \quad j = 1, \dots, m.$$

Define the symmetric bilinear form  $\langle \cdot, \cdot \rangle_m$  by

$$\langle \phi, \psi \rangle_m := \sum_{j=1}^m \mu_j \phi(c_j) \psi(c_j)$$

for any two polynomials  $\phi, \psi \in \mathcal{P}$ . This form is related to the form  $\langle \cdot, \cdot \rangle$  in an obvious way: instead of the zeros  $z_1, \dots, z_n$  and their multiplicities  $\nu_1, \dots, \nu_n$ , we now use the centres of gravity  $c_1, \dots, c_m$  and the weights  $\mu_1, \dots, \mu_m$  of the clusters. Let

$$\delta := \max_{1 \leq k \leq n} |\zeta_k|.$$



The following theorem tells us that  $\langle \cdot, \cdot \rangle_m$  approximates  $\langle \cdot, \cdot \rangle$  (and vice versa).

**THEOREM 3.1.** *Let  $\phi, \psi \in \mathcal{P}$ . Then  $\langle \phi, \psi \rangle = \langle \phi, \psi \rangle_m + \mathcal{O}(\delta^2)$ ,  $\delta \rightarrow 0$ .*

**PROOF.** The following holds:

$$\begin{aligned}
 \langle \phi, \psi \rangle &= \sum_{k=1}^n \nu_k \phi(z_k) \psi(z_k) \\
 &= \sum_{j=1}^m \sum_{k \in I_j} \nu_k \phi(c_j + \zeta_k) \psi(c_j + \zeta_k) \\
 &= \sum_{j=1}^m \sum_{k \in I_j} \nu_k \left( \phi(c_j) \psi(c_j) + \zeta_k [\phi(z) \psi(z)]'_{z=c_j} + \mathcal{O}(\zeta_k^2), \zeta_k \rightarrow 0 \right) \\
 &= \sum_{j=1}^m \mu_j \phi(c_j) \psi(c_j) + \underbrace{\sum_{j=1}^m \left( \sum_{k \in I_j} \nu_k \zeta_k \right) [\phi(z) \psi(z)]'_{z=c_j}}_{=0} + \sum_{k=1}^n \mathcal{O}(\zeta_k^2), \zeta_k \rightarrow 0 \\
 &= \sum_{j=1}^m \mu_j \phi(c_j) \psi(c_j) + \sum_{k=1}^n \mathcal{O}(\zeta_k^2), \zeta_k \rightarrow 0.
 \end{aligned}$$

This proves the theorem.  $\square$

Define the ordinary moments associated with  $\langle \cdot, \cdot \rangle_m$  as  $s_p^{(m)} := \langle 1, z^p \rangle_m$  for  $p = 0, 1, 2, \dots$ . Observe that  $s_0^{(m)} = s_0$  and  $s_1^{(m)} = s_1$ . Define the vectors  $\mathbf{s}, \mathbf{s}^{(m)} \in \mathbb{C}^{2N-1}$  as

$$\mathbf{s} := \begin{bmatrix} s_0 \\ s_1 \\ \vdots \\ s_{2N-2} \end{bmatrix} \quad \text{and} \quad \mathbf{s}^{(m)} := \begin{bmatrix} s_0^{(m)} \\ s_1^{(m)} \\ \vdots \\ s_{2N-2}^{(m)} \end{bmatrix}.$$

The entries of  $\mathbf{s}$  determine the Hankel matrix  $H_N$ . The previous theorem implies that

$$(3.1) \quad \frac{\|\mathbf{s} - \mathbf{s}^{(m)}\|_2}{\|\mathbf{s}\|_2} = \mathcal{O}(\delta^2), \delta \rightarrow 0.$$

Let  $H_k^{(m)}$  be the  $k \times k$  Hankel matrix

$$H_k^{(m)} := \left[ s_{p+q}^{(m)} \right]_{p,q=0}^{k-1}$$

for  $k = 1, 2, \dots$ .

**COROLLARY 3.2.** *Let  $k \geq 1$ . Then  $\det H_k = \det H_k^{(m)} + \mathcal{O}(\delta^2)$ ,  $\delta \rightarrow 0$ .*

**PROOF.** Let  $k$  be a positive integer. Then the previous theorem implies that

$$H_k = \left[ s_{p+q} \right]_{p,q=0}^{k-1} = \left[ s_{p+q}^{(m)} + \mathcal{O}(\delta^2), \delta \rightarrow 0 \right]_{p,q=0}^{k-1}.$$

The result follows by expanding the determinant of the matrix in the right-hand side.  $\square$

COROLLARY 3.3. *The matrix  $H_m$  is nonsingular if  $\delta \rightarrow 0$ . Let  $t > m$ . Then  $\det H_t = \mathcal{O}(\delta^2)$ ,  $\delta \rightarrow 0$ .*

PROOF. This follows from the previous corollary and the fact that  $H_m^{(m)}$  is nonsingular while  $H_t^{(m)}$  is singular for all integers  $t > m$  (cf. Theorem 2.1).  $\square$

The following theorem should be compared with Theorem 2.2.

THEOREM 3.4. *Let  $t$  be an integer  $\geq m$ . Then  $\varphi_t(c_j) = \mathcal{O}(\delta^2)$ ,  $\delta \rightarrow 0$  for  $j = 1, \dots, m$ . Also  $\langle z^p, \varphi_t(z) \rangle = \mathcal{O}(\delta^2)$ ,  $\delta \rightarrow 0$  for all  $p \geq t$ .*

PROOF. Let  $t \geq m$ . If  $t$  is a regular index, then

$$\langle z^p, \varphi_t(z) \rangle = 0, \quad p = 0, 1, \dots, t-1,$$

else

$$\langle z^p, \varphi_t(z) \rangle = 0, \quad p = 0, 1, \dots, r-1,$$

where  $r$  is the largest regular index less than  $t$ . Corollary 3.3 implies that  $r \geq m$ , and thus we may conclude that

$$\langle z^p, \varphi_t(z) \rangle = 0, \quad p = 0, 1, \dots, m-1.$$

Theorem 3.1 then implies that

$$\langle z^p, \varphi_t(z) \rangle_m = \mathcal{O}(\delta^2), \quad \delta \rightarrow 0, \quad p = 0, 1, \dots, m-1.$$

In matrix notation this can be written as

$$\begin{bmatrix} 1 & \cdots & 1 \\ c_1 & \cdots & c_m \\ \vdots & & \vdots \\ c_1^{m-1} & \cdots & c_m^{m-1} \end{bmatrix} \begin{bmatrix} \mu_1 & & \\ & \mu_2 & \\ & & \ddots \\ & & & \mu_m \end{bmatrix} \begin{bmatrix} \varphi_t(c_1) \\ \varphi_t(c_2) \\ \vdots \\ \varphi_t(c_m) \end{bmatrix} = \mathcal{O}(\delta^2), \quad \delta \rightarrow 0$$

where the right-hand side represents a vector in  $\mathbb{C}^m$  whose entries are  $\mathcal{O}(\delta^2)$ ,  $\delta \rightarrow 0$ . As the centres  $c_1, \dots, c_m$  are assumed to be mutually distinct and the weights  $\mu_1, \dots, \mu_m$  are different from zero, it follows that

$$\varphi_t(c_j) = \mathcal{O}(\delta^2), \quad \delta \rightarrow 0, \quad j = 1, \dots, m.$$

Theorem 3.1 then immediately implies that

$$\langle z^p, \varphi_t(z) \rangle = \mathcal{O}(\delta^2), \quad \delta \rightarrow 0$$

for all  $p \geq t$ .  $\square$

In other words, unless the FOP  $\varphi_t(z)$  has a very flat shape near its zeros, we are likely to find good approximations for the centres  $c_1, \dots, c_m$  among the zeros of  $\varphi_t(z)$  for all  $t \geq m$ . Note that

$$\begin{bmatrix} 1 & \cdots & 1 \\ c_1 & \cdots & c_m \\ \vdots & & \vdots \\ c_1^{m-1} & \cdots & c_m^{m-1} \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_m \end{bmatrix} = \begin{bmatrix} s_0 \\ s_1 \\ \vdots \\ s_{m-1} \end{bmatrix} + \mathcal{O}(\delta^2), \quad \delta \rightarrow 0.$$

It follows that approximations for the weights  $\mu_1, \dots, \mu_m$  can be obtained by solving a Vandermonde system.

#### 4 An accurate algorithm to compute zeros of FOPs.

We will now discuss an algorithm to compute zeros of FOPs. We consider the form  $\langle \cdot, \cdot \rangle$  as it can be evaluated explicitly via numerical integration. Our numerical experiments indicate that our algorithm gives very accurate results.

Define the matrices  $G_k$  and  $G_k^{(1)}$  as

$$G_k := [\langle \varphi_r, \varphi_s \rangle]_{r,s=0}^{k-1} \quad \text{and} \quad G_k^{(1)} := [\langle \varphi_r, \varphi_1 \varphi_s \rangle]_{r,s=0}^{k-1}$$

for  $k = 1, 2, \dots$ . Remember that the block orthogonality property of the FOPs implies that the Gram matrix  $G_n$  is block diagonal. The diagonal blocks are nonsingular, symmetric and zero above the main antidiagonal. If all the inner polynomials in a certain block are defined by using the standard monomial basis, then the corresponding diagonal block has Hankel structure. The matrix  $G_n^{(1)}$  is block tridiagonal. The diagonal blocks are symmetric and lower anti-Hessenberg (i.e., their entries are equal to zero along all the antidiagonals that lie above the main antidiagonal, except for the antidiagonal that precedes the main antidiagonal). Again, if all the inner polynomials in a certain block are defined by using the standard monomial basis, then the corresponding diagonal block is a Hankel matrix. The entries of the off-diagonal blocks are all equal to zero, except for the entry in the south-east corner. (For proofs and further details, we again refer to [13].)

The following theorem tells us that the zeros of a regular FOP can be calculated by solving a generalized eigenvalue problem. This will enable us to evaluate regular FOPs in their product representation, which is numerically very stable.

**THEOREM 4.1.** *Let  $t \geq 1$  be a regular index and let  $z_{t,1}, \dots, z_{t,t}$  be the zeros of the regular FOP  $\varphi_t$ . Then the eigenvalues of the pencil  $G_t^{(1)} - \lambda G_t$  are given by  $\varphi_1(z_{t,1}), \dots, \varphi_1(z_{t,t})$ . In other words, they are given by  $z_{t,1} - \mu, \dots, z_{t,t} - \mu$  where  $\mu = s_1/s_0$ .*

**PROOF.** Define the Hankel matrix  $H_t^<$  as  $H_t^< := [s_{1+k+l}]_{k,l=0}^{t-1}$ . We will first show that the zeros of  $\varphi_t$  are given by the eigenvalues of the pencil  $H_t^< - \lambda H_t$ . The zeros of  $\varphi_t$  are given by the eigenvalues of its companion matrix  $C_t$ . Let  $\lambda^*$  be an eigenvalue of  $C_t$  and  $x$  a corresponding eigenvector. As  $H_t$  is nonsingular, we may conclude that  $C_t x = \lambda^* x \Leftrightarrow H_t C_t x = \lambda^* H_t x$ . Using (2.3) one can easily verify that  $H_t C_t = H_t^<$ .

Let  $A_t$  be the unit upper triangular matrix that contains the coefficients of  $\varphi_0, \varphi_1, \dots, \varphi_{t-1}$ . Then  $G_t$  can be factorized as  $G_t = A_t^T H_t A_t$ . As  $\varphi_1(z) = z - \mu$  where  $\mu = s_1/s_0$ , the matrix  $G_t^{(1)}$  is given by  $[\langle \varphi_r, z \varphi_s \rangle]_{r,s=0}^{t-1} - \mu G_t$ . The matrix  $[\langle \varphi_r, z \varphi_s \rangle]_{r,s=0}^{t-1}$  can be written as  $A_t^T H_t^< A_t$  and thus  $G_t^{(1)} = A_t^T (H_t^< - \mu H_t) A_t$ . Now let  $\lambda^*$  be an eigenvalue of the pencil  $H_t^< - \lambda H_t$  and  $x$  a corresponding

eigenvector. Then

$$\begin{aligned} H_t^< x &= \lambda^* H_t x \\ \Leftrightarrow (H_t^< - \mu H_t) x &= (\lambda^* - \mu) H_t x \\ \Leftrightarrow A_t^T (H_t^< - \mu H_t) A_t y &= \varphi_1(\lambda^*) A_t^T H_t A_t y \quad \text{if } y := A_t^{-1} x \\ \Leftrightarrow G_t^{(1)} y &= \varphi_1(\lambda^*) G_t y. \end{aligned}$$

This proves the theorem.  $\square$

Regular FOPs are characterized by the fact that the determinant of a Hankel matrix is different from zero, while inner polynomials correspond to singular Hankel matrices. To decide whether  $\varphi_t(z)$  should be defined as a regular FOP or as an inner polynomial, one could therefore calculate the determinant of  $H_t$  and check if it is equal to zero. However, from a numerical point of view such a test “is equal to zero” does not make sense. Because of rounding errors (both in the evaluation of  $\langle \cdot, \cdot \rangle$  and the calculation of the determinant) we would encounter only regular FOPs. Strictly speaking one could say that inner polynomials are not needed in numerical calculations. However, the opposite is true! Let us agree to call a regular FOP *well-conditioned* if its corresponding Yule–Walker system (2.3) is well-conditioned, and *ill-conditioned* otherwise. To obtain a numerically stable algorithm, it is crucial to generate only well-conditioned regular FOPs and to replace ill-conditioned regular FOPs by inner polynomials. Stable look-ahead solvers for linear systems of equations that have Hankel structure are based on this principle [11, 15, 19]. In this approach the diagonal blocks in  $G_n$  are taken (slightly) larger than strictly necessary to avoid ill-conditioned blocks. A disadvantage is that part of the structure of  $G_n$  and  $G_n^{(1)}$  gets lost, i.e., there will be some additional fill-in.

Our algorithm for computing zeros of analytic functions proceeds by calculating the polynomials  $\varphi_0(z), \varphi_1(z), \dots, \varphi_n(z)$  in their product representation, starting with  $\varphi_0(z) \leftarrow 1$  and  $\varphi_1(z) \leftarrow z - \mu$ . At each step we ask ourselves whether it is numerically feasible to generate the next polynomial in the sequence as a regular FOP. As the reader will see, there are several ways to find an answer to this question.

The value of  $n$  is determined as follows. Suppose that the algorithm has just generated a (well-conditioned) regular FOP  $\varphi_r(z)$ . To check whether  $n = r$ , we scan the sequence

$$\left( |\langle (z - \mu)^\tau \varphi_r(z), \varphi_r(z) \rangle| \right)_{\tau=0}^{N-1-r};$$

cf. Theorem 2.2. If all the elements are “sufficiently small”, then we conclude that indeed  $n = r$  and we stop.

The form  $\langle \cdot, \cdot \rangle$  is evaluated via numerical integration along  $\gamma$ , i.e., it is approximated by a quadrature sum. We assume that this sum is calculated in the standard way, by adding the terms one by one, in other words, by forming a sequence of partial sums. We ask the quadrature algorithm not only for an approximation of the integral, say **result**, but also for the modulus of the partial

sum that has the largest modulus, say **maxpsum**. Then

$$\log_{10} \frac{\mathbf{maxpsum}}{|\mathbf{result}|}$$

is an estimate for the loss of precision. This information will turn out to be extremely useful, for example in the stopping criterion.

These considerations lead to the following algorithm.

ALGORITHM 1.

**input**  $\langle \cdot, \cdot \rangle, \epsilon_{\text{stop}}$   
**output**  $n, \text{zeros}$   
**comment**  $\text{zeros} = \{z_1, \dots, z_n\}$ . We assume that  $\epsilon_{\text{stop}} > 0$ .  
 $N \leftarrow \langle 1, 1 \rangle$

**if**  $N == 0$  **then**  
      $n \leftarrow 0$ ;  $\text{zeros} \leftarrow \emptyset$ ; **stop**  
**else**  
      $\varphi_0(z) \leftarrow 1$   
      $\mu \leftarrow \langle z, 1 \rangle / N$ ;  $\varphi_1(z) \leftarrow z - \mu$   
      $r \leftarrow 1$ ;  $t \leftarrow 0$   
     **while**  $r + t < N$  **do**  
         regular  $\leftarrow$  it is numerically feasible to generate  $\varphi_{r+t+1}(z)$  as  
         a regular FOP      ... [1]  
         **if** regular **then**  
             generate  $\varphi_{r+t+1}(z)$  as a regular FOP      ... [2]  
              $r \leftarrow r + t + 1$ ;  $t \leftarrow 0$   
             allsmall  $\leftarrow$  **true**;  $\tau \leftarrow 0$   
             **while** allsmall **and**  $(r + \tau < N)$  **do**  
                  $[\text{ip}, \text{maxpsum}] \leftarrow \langle (z - \mu)^\tau \varphi_r(z), \varphi_r(z) \rangle$       ... [3]  
                  $\text{ip} \leftarrow |\text{ip}|$   
                 allsmall  $\leftarrow (\text{ip} / \text{maxpsum} < \epsilon_{\text{stop}})$       ... [4]  
                  $\tau \leftarrow \tau + 1$   
             **end while**  
             **if** allsmall **then**  
                  $n \leftarrow r$ ;  $\text{zeros} \leftarrow \text{roots}(\varphi_r)$ ; **stop**  
             **end if**  
         **else**  
             generate  $\varphi_{r+t+1}(z)$  as an inner polynomial      ... [5]  
              $t \leftarrow t + 1$   
         **end if**  
     **end while**  
      $n \leftarrow N$ ;  $\text{zeros} \leftarrow \text{roots}(\varphi_N)$ ; **stop**  
**end if**

COMMENTS.

1. Statement [1] is crucial. But how does one decide that it is numerically feasible to generate the next polynomial in the sequence

$$\varphi_0(z), \varphi_1(z), \dots, \varphi_n(z)$$

as a regular FOP? Suppose that the algorithm has just generated a regular FOP  $\varphi_r(z)$ . By using an explicit determinant expression for regular FOPs, one can show that  $\langle \varphi_r(z), \varphi_r(z) \rangle = \det H_{r+1} / \det H_r$ . Therefore, if  $r$  is a regular index, then  $r+1$  is a regular index if and only if  $\langle \varphi_r(z), \varphi_r(z) \rangle \neq 0$ . This suggests the following criterion: if  $|\langle \varphi_r(z), \varphi_r(z) \rangle| / \text{maxsum} < \epsilon_{\text{regular}}$ , where  $\epsilon_{\text{regular}}$  is some small threshold given by the user, then define  $\varphi_{r+1}(z)$  as an inner polynomial, else define it as a regular FOP. However, as we will illustrate in the next section, it is very difficult to choose an appropriate value of  $\epsilon_{\text{regular}}$ .

We prefer to use the following criterion: act as if the next polynomial in the sequence, say  $\varphi_t(z)$ , is defined as a regular FOP, i.e., compute its zeros by computing the eigenvalues of the pencil  $G_t^{(1)} - \lambda G_t$  and then check if these zeros lie sufficiently close to the interior of  $\gamma$ . If so, then define  $\varphi_t(z)$  as a regular FOP, else define it as an inner polynomial. The idea behind this strategy is the following. If the matrix  $G_t$  is singular, in which case also the matrix  $H_t$  is singular of course, then the pencil  $G_t^{(1)} - \lambda G_t$  has either a number of eigenvalues at infinity or a number of eigenvalues that may assume arbitrary values. Indeed, by using the structure of the matrices  $G_t^{(1)}$  and  $G_t$  one can easily prove the following result, which complements Theorem 4.1.

**THEOREM 4.2.** *Let  $t \geq 1$  be an integer, let  $r$  be the largest regular index less than or equal to  $t$ , and let  $r'$  be the smallest regular index greater than  $t$ . (Define  $r' := +\infty$  if  $t \geq n$ .) Then the eigenvalues of the pencil  $G_t^{(1)} - \lambda G_t$  are given by the eigenvalues of the pencil  $G_r^{(1)} - \lambda G_r$  and  $t-r$  eigenvalues that may assume arbitrary values if  $t < r' - 1$  or  $t-r$  eigenvalues  $\lambda = \infty$  if  $t = r' - 1$ .*

Each of these  $t-r$  indeterminate eigenvalues corresponds to two corresponding zeros on the diagonals of the generalized Schur decomposition of  $G_t^{(1)}$  and  $G_t$ . When actually calculated, these diagonal entries are different from zero because of roundoff errors. The quotient of two such corresponding diagonal entries is a spurious eigenvalue. Our strategy is based on the assumption that, if the matrix  $H_t$ , and thus also the matrix  $G_t$ , is nearly singular, then the computed eigenvalues of the pencil  $G_t^{(1)} - \lambda G_t$  that correspond to the eigenvalues that lie at infinity or that may assume arbitrary values, lie far away from the interior of  $\gamma$ .

The reader may object that our criterion is too strict. Indeed, the zeros of the regular FOPs of degree  $< n$  need not lie close to the interior of  $\gamma$ , except if the form  $\langle \cdot, \cdot \rangle$  is a true (positive definite) inner product, in which case the zeros of the regular FOPs lie in the convex hull of  $\{z_1, \dots, z_n\}$ . (This follows from a general result on orthogonal polynomials. See, e.g., Van Assche [50].) Thus it may very well be the case that some of the computed zeros of a well-conditioned regular FOP lie far away from  $\gamma$ , in which case our algorithm decides

to define this polynomial as an inner polynomial. In other words, our algorithm may define more inner polynomials than strictly necessary. However, we have done a lot of numerical tests, and have found that our strategy leads to very accurate results. Also, compared to the criterion based on inner products of the type  $\langle \varphi_r(z), \varphi_r(z) \rangle$ , another advantage is that the user doesn't have to supply a threshold such as  $\epsilon_{\text{regular}}$ .

2. Statement [2] means: define  $\varphi_{r+t+1}(z)$  as  $\varphi_{r+t+1}(z) \leftarrow \prod_{j=1}^{r+t+1} (z - \alpha_j)$ . The zeros  $\alpha_j$  are computed as  $\alpha_j = \mu + \lambda_j$ ,  $j = 1, \dots, r+t+1$ , where  $\lambda_1, \dots, \lambda_{r+t+1}$  are the eigenvalues of the pencil  $G_{r+t+1}^{(1)} - \lambda G_{r+t+1}$ ; cf. Theorem 4.1.
3. In statement [3] we use the inner product  $\langle (z - \mu)^\tau \varphi_r(z), \varphi_r(z) \rangle$  and not  $\langle z^\tau \varphi_r(z), \varphi_r(z) \rangle$  as it is likely that the former leads to more accurate results than the latter. If  $\tau \leq r$ , then one may also use  $\langle \varphi_\tau(z) \varphi_r(z), \varphi_r(z) \rangle$ . In general, if  $\tau = \alpha r + \beta$ , where  $\alpha, \beta \in \mathbb{N}$  with  $\beta < r$ , then one may use  $\langle [\varphi_r(z)]^{\alpha+1} \varphi_\beta(z), \varphi_r(z) \rangle$ .
4. Observe that in statement [4] we do not compare  $\text{ip}$  with  $\epsilon_{\text{stop}}$  but take into account the loss of precision as estimated by the quadrature algorithm. We have found this heuristic to be very reliable.
5. In statement [5] one may define  $\varphi_{r+t+1}(z)$  as  $\varphi_{r+t+1}(z) \leftarrow (z - \mu) \varphi_{r+t}(z)$  or  $\varphi_{r+t+1}(z) \leftarrow \varphi_{t+1}(z) \varphi_r(z)$ . Both versions are to be preferred to the "classical"  $\varphi_{r+t+1}(z) \leftarrow z^{t+1} \varphi_r(z)$ .
6. Instead of computing  $\mu$ , the arithmetic mean of the zeros, as  $\mu \leftarrow \langle z, 1 \rangle / N$ , one can also use the following formula, which may give a more accurate result:  $\mu \leftarrow w + \langle z - w, 1 \rangle / N$ , where  $w$  is a point inside  $\gamma$ , preferably near the centre of the interior of  $\gamma$ .
7. As we represent our FOPs by using the product representation,  $\varphi(z) = \prod_{\alpha \in \varphi^{-1}(0)} (z - \alpha)$ , the function  $\text{roots}(\cdot)$  is obviously *not* a function that calculates the zeros of a polynomial from its coefficients in the standard monomial basis.

What happens if we apply our algorithm in case the zeros of  $f$  can be grouped into clusters? The second part of Theorem 3.4 implies that our algorithm stops at  $r = m$  if  $\delta$ , the maximal size of the clusters, is sufficiently small. It returns the zeros of the FOP  $\varphi_m(z)$  that is associated with  $\langle \cdot, \cdot \rangle$ . Theorem 3.1 and the fact that the  $m$ th degree FOP with respect to  $\langle \cdot, \cdot \rangle_m$  is given by  $\prod_{j=1}^m (z - c_j)$  imply that we can use these zeros as approximations  $\hat{c}_1, \dots, \hat{c}_m$  for the centres of the clusters. (This also follows from the first part of Theorem 3.4, of course.) By solving the Vandermonde system

$$\begin{bmatrix} 1 & \cdots & 1 \\ \hat{c}_1 & \cdots & \hat{c}_m \\ \vdots & & \vdots \\ \hat{c}_1^{m-1} & \cdots & \hat{c}_m^{m-1} \end{bmatrix} \begin{bmatrix} \hat{\mu}_1 \\ \hat{\mu}_2 \\ \vdots \\ \hat{\mu}_m \end{bmatrix} = \begin{bmatrix} s_0 \\ s_1 \\ \vdots \\ s_{m-1} \end{bmatrix}$$

we obtain approximations  $\hat{\mu}_1, \dots, \hat{\mu}_m$  for the weights of the clusters. These should be close to integers. We can check this, to verify that we have indeed determined the correct value of  $m$ . We can also calculate (approximations for) the ordinary moments associated with  $\langle \cdot, \cdot \rangle_m$  and verify if (3.1) is satisfied.

## 5 Numerical examples.

We have implemented our algorithm in MATLAB. The m-files are available from the authors. In the following examples, the computations have been done via MATLAB 5 (with floating-point relative accuracy  $\approx 2.2204 \cdot 10^{-16}$ ).

We have considered the case that  $\gamma$  is a circle. The following integration algorithm is used to approximate the form  $\langle \cdot, \cdot \rangle$ . Suppose that  $\gamma$  is the circle with centre  $c$  and radius  $\rho$ . Then

$$(5.1) \quad \langle \phi, \psi \rangle = \rho \int_0^1 \phi(c + \rho e^{2\pi i \theta}) \psi(c + \rho e^{2\pi i \theta}) \frac{f'(c + \rho e^{2\pi i \theta})}{f(c + \rho e^{2\pi i \theta})} e^{2\pi i \theta} d\theta.$$

Since this is the integral of a periodic function over a complete period, the trapezoidal rule is an appropriate quadrature rule. If  $F : [0, 1] \rightarrow \mathbb{C}$  is the integrand in the right-hand side of (5.1), then the  $q$ -point trapezoidal rule approximation to  $\langle \phi, \psi \rangle$  is given by

$$\langle \phi, \psi \rangle = \int_0^1 F(\theta) d\theta \approx \frac{1}{q} \sum_{k=0}^{q-1} F(k/q) =: T_q.$$

The double prime indicates that the first and the last term of the sum are to be multiplied by  $1/2$ . As  $F$  is periodic with period one, we may rewrite  $T_q$  as

$$T_q = \frac{1}{q} \sum_{k=0}^{q-1} F(k/q).$$

This shows that  $T_q$  indeed depends on  $q$  (and not  $q+1$ ) points. As

$$T_{2q} = \frac{1}{2} T_q + T_{q \rightarrow 2q}$$

where

$$T_{q \rightarrow 2q} := \frac{1}{2q} \sum_{k=0}^{q-1} F\left(\frac{2k+1}{2q}\right),$$

successive doubling of  $q$  enables us in each step to reuse the integrand values needed in the previous step. In the following examples we started with  $q = 16$  and continued doubling  $q$  until  $|T_{2q} - T_q|$  was sufficiently small. More precisely, if  $S_q$  and  $S_{q \rightarrow 2q}$  denote the modulus of the partial sum of  $qT_q$  respectively  $2qT_{q \rightarrow 2q}$  that has the largest modulus, then our stopping criterion is given by  $|T_{2q} - T_q| \leq S_{2q} 10^{-14}$ , where  $S_{2q} := \max\{S_q, S_{q \rightarrow 2q}\}/(2q)$ .

Lyness and Delves [40] studied the asymptotic behaviour of the quadrature error. They showed that the modulus of the error made by the  $q$ -point trapezoidal rule is asymptotically  $\mathcal{O}(A^q)$  where  $0 \leq A < 1$ . More precisely,

$$A := \max \left\{ \frac{|z_I|}{\rho}, \frac{\rho}{|z_E|}, \frac{\rho}{\rho_s} \right\}$$



where  $z_I$  is the zero of  $f$  that lies closest to  $\gamma$  and in the interior of  $\gamma$ ,  $z_E$  is the zero of  $f$  that lies closest to  $\gamma$  and in the exterior of  $\gamma$ , and  $\rho_s$  is the distance between  $c$  and the nearest singularity of  $f$ .

EXAMPLE 5.1. Our first example illustrates the importance of shifting the origin in the complex plane to the arithmetic mean of the zeros. It also compares the two strategies that we have proposed to decide whether it is numerically feasible to generate the next polynomial in the sequence  $\varphi_0(z), \varphi_1(z), \dots, \varphi_n(z)$  as a regular FOP. We will see that it is indeed better to act as if the polynomial is a regular FOP, i.e., to compute its zeros by solving the generalized eigenvalue problem of Theorem 4.1, and then to check if these zeros lie sufficiently close to the interior of  $\gamma$ . If so, the polynomial is indeed defined as a regular FOP, else it is defined as an inner polynomial.

Suppose that  $n = 3$ ,  $z_1 = \epsilon$ ,  $z_2 = \sqrt{3} + i$ ,  $z_3 = \sqrt{3} - i$ , and  $\nu_1 = \nu_2 = \nu_3 = 1$ . That is, suppose that  $f(z) = (z - \epsilon)[(z - \sqrt{3})^2 + 1]$ . If  $\epsilon = 0$ , then the Hankel matrix  $H_2$  is exactly singular, i.e.,  $\varphi_2(z)$  has to be defined as an inner polynomial. We set  $\epsilon = 10^{-2}$ . Suppose that  $\gamma = \{z \in \mathbb{C} : |z| = 3\}$ . In the quadrature algorithm, we have evaluated the logarithmic derivative  $f'(z)/f(z)$  of  $f(z)$  via the formula

$$(5.2) \quad \frac{f'(z)}{f(z)} = \sum_{k=1}^n \frac{\nu_k}{z - z_k}.$$

We have taken  $\epsilon_{\text{stop}} = 10^{-18}$ . Our algorithm proceeds as follows. The total number of zeros  $N$  is equal to 3. The polynomial  $\varphi_0(z)$  is of course defined as a regular FOP,  $\varphi_0(z) \leftarrow 1$ . The arithmetic mean  $\mu$  is approximated by

$$1.158033871706381 \text{ e}+00 + \text{i} \cdot 1.526556658859590 \text{ e}-16.$$

The polynomial  $\varphi_1(z)$  is also defined as a regular FOP,  $\varphi_1(z) \leftarrow z - \mu$ . The inner product  $\langle \varphi_1(z), \varphi_1(z) \rangle$  is equal to 0.02303. To take into account the loss of precision, we divide by  $S_{2q}$  to obtain its scaled counterpart, which is equal to 0.01231. Is this that small that we should define  $\varphi_2(z)$  as an inner polynomial? It seems not, and we decide to define  $\varphi_2(z)$  as a regular FOP. Its zeros are approximated by

$$\begin{aligned} &1.158072473165547 \text{ e}+00 + \text{i} \cdot 1.513258564730580 \text{ e}-16 \\ &2.000049565733722 \text{ e}+02 + \text{i} \cdot 3.487413819470906 \text{ e}-12 \end{aligned}$$

Note how large the second zero is! The inner product  $\langle \varphi_2(z), \varphi_2(z) \rangle$  is equal to 910.504. Its scaled counterpart is 0.01214. We decide to define  $\varphi_3(z)$  as a regular FOP. Its zeros are approximated by

$$\begin{aligned} &1.732050807571817 \text{ e}+00 - \text{i} \cdot 1.000000000004209 \text{ e}+00 \\ &1.000000000661760 \text{ e}-02 + \text{i} \cdot 3.913802997325630 \text{ e}-12 \\ &1.732050807566777 \text{ e}+00 + \text{i} \cdot 1.000000000002807 \text{ e}+00 \end{aligned}$$

As  $N = 3$ , we may stop. The relative errors of the approximations for the zeros of  $f$  are  $\mathcal{O}(10^{-12})$ , except for the zero that approximates  $z_1 = \epsilon$ , which has a relative error of  $\mathcal{O}(10^{-10})$ . The absolute errors are  $\mathcal{O}(10^{-12})$ . The relative errors of the approximations for the multiplicities are  $\mathcal{O}(10^{-11})$ .

As one of the zeros of  $\varphi_2(z)$  lies far away from the interior of  $\gamma$ , we should decide to define  $\varphi_2(z)$  as an inner polynomial. Surprisingly, this does not improve

the accuracy of the results. However, let us see what happens if we first shift the origin to  $\mu$ , or, equivalently, if we consider the circle  $\gamma = \{z \in \mathbb{C} : |z - \mu| = 2\}$ . Note that we change both the centre and the radius of  $\gamma$ . By defining  $\varphi_2(z)$  as a regular FOP, the accuracy of the results does not improve. However, if we define  $\varphi_2(z)$  as an inner polynomial, the relative errors of the approximations for the zeros of  $f$  are  $\mathcal{O}(10^{-16})$ , except for the zero that approximates  $z_1 = \epsilon$ , which has a relative error of  $\mathcal{O}(10^{-14})$ . The absolute errors are  $\mathcal{O}(10^{-16})$ . The relative errors of the approximations for the multiplicities are  $\mathcal{O}(10^{-16})$ . In other words, the results are indeed much better.  $\diamond$

EXAMPLE 5.2. Let  $f(z) = e^{3z} + 2z \cos z - 1$  and  $\gamma = \{z \in \mathbb{C} : |z| = 2\}$ . We set  $\epsilon_{\text{stop}} = 10^{-18}$ . Our algorithm finds that  $N = 4$ . It defines  $\varphi_0(z)$  and  $\varphi_1(z)$  as regular FOPs. From the eigenvalues of the pencil  $G_2^{(1)} - \lambda G_2$  it concludes that  $\varphi_2(z)$  would have a zero of modulus  $\approx 43$  in case  $\varphi_2(z)$  is defined as a regular FOP. Thus the algorithm decides to define  $\varphi_2(z)$  as an inner polynomial. The polynomials  $\varphi_3(z)$  and  $\varphi_4(z)$  are defined as regular FOPs. Our algorithm concludes that  $n = 4$ . The computed approximations for the zeros of  $f$  are given by

$$\begin{aligned} & -1.844233953262213 \text{ e}+00 - \text{i} \cdot 1.106288924192872 \text{ e}-16 \\ & 5.308949302929297 \text{ e}-01 + \text{i} \cdot 1.331791876751121 \text{ e}+00 \\ & 5.308949302929303 \text{ e}-01 - \text{i} \cdot 1.331791876751121 \text{ e}+00 \\ & -5.412337245047638 \text{ e}-15 + \text{i} \cdot 3.762630283199076 \text{ e}-16 \end{aligned}$$

The corresponding approximations for the multiplicities are

$$\begin{aligned} & 1.0000000000000001 \text{ e}+00 + \text{i} \cdot 9.279422312879846 \text{ e}-17 \\ & 1.0000000000000001 \text{ e}+00 - \text{i} \cdot 2.415808667423342 \text{ e}-15 \\ & 1.0000000000000001 \text{ e}+00 + \text{i} \cdot 1.187431378902999 \text{ e}-15 \\ & 9.999999999999974 \text{ e}-01 + \text{i} \cdot 1.142195850770565 \text{ e}-15 \end{aligned}$$

By refining the approximations for the zeros of  $f$  via Newton's method, we find that they have a relative error of  $\mathcal{O}(10^{-16})$ , except for the approximation of  $z_4 = 0$ , which has an absolute error of  $\mathcal{O}(10^{-15})$ . If  $\varphi_2(z)$  is defined as a regular FOP, the errors are  $\mathcal{O}(10^{-13})$ .  $\diamond$

Stewart's perturbation theory for the generalized eigenvalue problem [48] allows us to make a sensitivity analysis. The main result from his first order perturbation theory for simple eigenvalues tells us the following. If  $\lambda$  is a simple eigenvalue of the pencil  $G_t^{(1)} - \lambda G_t$  and  $\lambda_\epsilon$  is the corresponding eigenvalue of a perturbed pencil  $\tilde{G}_t^{(1)} - \lambda \tilde{G}_t$  with  $\|G_t^{(1)} - \tilde{G}_t^{(1)}\|_2 \approx \|G_t - \tilde{G}_t\|_2 \approx \epsilon$ , then

$$\frac{|\lambda - \lambda_\epsilon|}{\sqrt{1 + |\lambda|^2} \sqrt{1 + |\lambda_\epsilon|^2}} \leq \frac{\epsilon}{\sqrt{|y^H G_t^{(1)} x|^2 + |y^H G_t x|^2}} + \mathcal{O}(\epsilon^2) =: \kappa(\lambda, x, y) \epsilon + \mathcal{O}(\epsilon^2)$$

where  $x$  and  $y$  are the right and left eigenvectors corresponding to  $\lambda$ ,

$$G_t^{(1)} x = \lambda G_t x \quad \text{and} \quad y^H G_t^{(1)} = \lambda y^H G_t,$$

normalized such that  $\|x\|_2 = \|y\|_2 = 1$ . Let us call  $\kappa(\lambda, x, y)$  the *sensitivity factor* of  $\lambda$ .

EXAMPLE 5.3. Suppose that  $n = 10$ ,

$$\begin{aligned} z_1 &= -1, \\ z_2 &= 4, \quad z_3 = 4 + \delta(1 + i), \\ z_4 &= 3i, \quad z_5 = 3i + \delta(10 + 5i), \quad z_6 = 3i + \delta(-3 + 4i), \\ z_7 &= c + \delta(-1 + 2i), \quad z_8 = c + \delta(1 + 5i), \quad z_9 = c + \delta(1 + i), \quad z_{10} = c + \delta(-2 - 2i), \end{aligned}$$

where  $c = -3 + 3i$  and  $\delta = 10^{-4}$ . Suppose that  $\nu_1 = \dots = \nu_{10} = 1$ . Let  $f(z)$  be the polynomial that has  $z_1, \dots, z_{10}$  as simple zeros,  $f(z) = \prod_{k=1}^{10} (z - z_k)$ , and let  $\gamma = \{z \in \mathbb{C} : |z| = 5\}$ . Note that  $f$  has four clusters of zeros,  $m = 4$ , of weight 1, 2, 3 and 4, respectively. We have evaluated the logarithmic derivative of  $f(z)$  again via formula (5.2). As we do not want the algorithm to stop as soon as it has found the approximations for the centres of the clusters, we set  $\epsilon_{\text{stop}}$  to a rather small value,  $\epsilon_{\text{stop}} = 10^{-18}$ . Our algorithm gives the following results. The total number of zeros is equal to 10. The polynomials  $\varphi_0(z)$  and  $\varphi_1(z)$  are defined as regular FOPs. The computed eigenvalues of the pencil  $G_2^{(1)} - \lambda G_2$  lead to zeros that lie inside  $\gamma$ , and thus the polynomial  $\varphi_2(z)$  is defined as a regular FOP. The sensitivity factors of the eigenvalues are equal to

3.773765047881042 e-02  
9.474367914383353 e-03

The solution of the Vandermonde system that corresponds to the computed approximations for the zeros of  $\varphi_2(z)$  is given by

7.050510110011339 e+00 - i. 2.941868186347331 e-01  
2.949489889988664 e+00 + i. 2.941868186347336 e-01

The algorithm computes  $\langle \varphi_2(z), \varphi_2(z) \rangle$ . In step [4] it compares

9.853283227226082 e-01

with  $\epsilon_{\text{stop}}$  and sets **allsmall**  $\leftarrow$  **false**. The polynomial  $\varphi_3(z)$  is defined as a regular FOP. The sensitivity factors of the eigenvalues are equal to

5.221220847969128 e-03  
5.691118489264763 e-02  
9.440130958262826 e-03

The solution of the Vandermonde system that corresponds to the computed approximations for the zeros of  $\varphi_3(z)$  is given by

1.751672615711886 e+00 + i. 2.391215777821213 e-02  
5.761474298212248 e+00 + i. 1.418723240124189 e-01  
2.486853086075870 e+00 - i. 1.657844817906288 e-01

The algorithm computes  $\langle \varphi_3(z), \varphi_3(z) \rangle$ . In step [4] it compares

9.560608076054004 e-02

with  $\epsilon_{\text{stop}}$  and sets **allsmall**  $\leftarrow$  **false**. The polynomial  $\varphi_4(z)$  is defined as a regular FOP. The sensitivity factors of the eigenvalues are equal to

5.997609084856927 e-03  
1.632514220886046 e-02  
2.198497992029341 e-02  
4.923118129602375 e-03

The solution of the Vandermonde system that corresponds to the computed approximations for the zeros of  $\varphi_4(z)$  is given by

```
1.999999998746185 e+00 - i·1.260274035855700e-08
4.000000263704519 e+00 - i·1.428161173014011e-07
2.999999744825852 e+00 + i·2.437473983804684e-07
9.999999927234455 e-01 - i·8.832854039794525e-08
```

Observe that these “multiplicities” (actually, they are the weights of the clusters) are at a distance of  $\mathcal{O}(10^{-8}) = \mathcal{O}(\delta^2)$  to integers. This is a first indication of the fact that  $m = 4$ . The algorithm computes  $\langle \varphi_4(z), \varphi_4(z) \rangle$ . In step [4] it compares

```
4.690246227384357 e-09
```

with  $\epsilon_{\text{stop}}$ . As we have given  $\epsilon_{\text{stop}}$  a very small value,  $\epsilon_{\text{stop}} = 10^{-18}$ , the algorithm sets `allsmall`  $\leftarrow$  **false** and continues. It defines the polynomial  $\varphi_5(z)$  as a regular FOP. The sensitivity factors of the eigenvalues are equal to

```
1.859124283121160 e+02
5.997664519337178 e-03
1.632900762729697 e-02
4.930300860498997 e-03
2.201051194773096 e-02
```

Observe that one of the eigenvalues is much more sensitive than the others. The solution of the Vandermonde system that corresponds to the computed approximations for the zeros of  $\varphi_5(z)$  is given by

```
-4.028597029147776 e-08 - i·1.692141150994100 e-07
2.000000007231077 e+00 - i·8.308500282558695 e-09
4.000000510541128 e+00 - i·2.651932814361151 e-07
1.000000192451618 e+00 - i·8.614030726982243 e-08
2.999999330062146 e+00 + i·5.288562050783197 e-07
```

Observe that the component that corresponds to the spurious eigenvalue is of size  $\mathcal{O}(10^{-8})$  whereas the other components are close to integers. This enables us to deduce the presence of spurious eigenvalues without computing the sensitivity factors. The algorithm computes  $\langle \varphi_5(z), \varphi_5(z) \rangle$ . In step [4] it compares

```
3.544154048709335 e-10
```

with  $\epsilon_{\text{stop}}$  and sets `allsmall`  $\leftarrow$  **false**. The polynomial  $\varphi_6(z)$  is defined as a regular FOP. The sensitivity factors of the eigenvalues are equal to

```
1.617034259272598 e+02
1.220513781306331 e+01
5.997660002011287 e-03
1.633002045462034 e-02
4.929209225445833 e-03
2.200115675959345 e-02
```

The solution of the Vandermonde system that corresponds to the computed approximations for the zeros of  $\varphi_6(z)$  is given by

```
-1.109848867030467 e-07 - i·3.954458493988473 e-08
5.068148299128812 e-12 - i·6.003619075497872 e-13
2.000000007044562 e+00 - i·9.500131431430428 e-09
```

4.000000455291260 e+00 - i· 3.668720544062149 e-07  
 1.000000170753607 e+00 - i· 1.185972639722307 e-07  
 2.999999477890392 e+00 + i· 5.345146309560269 e-07

The algorithm computes  $\langle \varphi_6(z), \varphi_6(z) \rangle$ . In step [4] it compares

1.189012222825083 e-09

with  $\epsilon_{\text{stop}}$  and sets  $\text{allsmall} \leftarrow \mathbf{false}$ . And so on. The algorithm defines  $\varphi_7(z)$  as a regular FOP,  $\varphi_8(z)$  as an inner polynomial,  $\varphi_9(z)$  as a regular FOP, and finally  $\varphi_{10}(z)$  as an inner polynomial. The computed approximations for the zeros of  $f$  are given by

-9.99999999599027 e-01 + i· 1.827471507453993 e-11  
 4.000066008247924 e+00 + i· 5.012986226260452 e-05  
 4.000047628710319 e+00 - i· 2.622572865154105 e-04  
 -5.868580001572310 e-05 + i· 3.000189455314092 e+00  
 1.147726665656712 e-03 + i· 3.000342939244789 e+00  
 -3.000001882960546 e+00 + i· 3.000324337949845 e+00  
 -3.000258740452308 e+00 + i· 2.999857462002724 e+00  
 -9.085700394437596 e-01 + i· 1.866967988946470 e+00  
 -8.509233181288104 e-01 - i· 6.684482268880147 e+00  
 -4.999300000000002 e-01 + i· 2.100160000000000 e+00

The relative errors of the approximations for the zeros that belong to the clusters of weight 1 and 2 are  $\mathcal{O}(10^{-11})$  and  $\mathcal{O}(10^{-5})$ , resp. For the other zeros, the relative errors are at least  $\mathcal{O}(10^{-3})$ .

If we set  $\epsilon_{\text{stop}} = 10^{-6}$ , then our algorithm stops at the polynomial of degree 4. We obtain the following approximations for the centres of the clusters:

-9.999999564181510 e-01 - i· 5.152524762408461 e-08  
 4.000050001653271 e+00 + i· 5.000694739720757 e-05  
 2.335838430156945 e-04 + i· 3.000299920075392 e+00  
 -3.000024926663507 e+00 + i· 3.000149946356108 e+00

Let us now focus on the separate clusters. We have considered the circles whose centre is the computed approximation for the centre of a cluster and whose radius is equal to 0.1. The relative errors of the approximations for the zeros that we obtain are  $\mathcal{O}(10^{-16})$ ,  $\mathcal{O}(10^{-12})$ ,  $\mathcal{O}(10^{-10})$  and  $\mathcal{O}(10^{-6})$  for the clusters of weight 1, 2, 3 and 4, respectively. If we consider a circle whose centre is the computed approximation for the centre of the cluster of weight 4 and whose radius is equal to  $10^{-3}$ , then the relative errors of the approximations that we obtain for the zeros that lie in this cluster are  $\mathcal{O}(10^{-16})$ . Apparently, the smaller the radius of the circle is, the more accurate the computed approximations for the zeros are. This can be explained by the fact that the quadrature method gives more accurate approximations for the integrals (in other words, for the data from which approximations for the zeros are computed).  $\diamond$

More numerical examples will be given in Section 7.

## 6 Rational interpolation at roots of unity.

We will now approach our problem of computing all the zeros of  $f$  that lie inside  $\gamma$  in a different way, based on rational interpolation at roots of unity. Interesting connections exist with the techniques presented in Sections 2, 3 and 4. We will show how the new approach complements the previous one.

Let  $K$  be a positive integer and let  $t_1, \dots, t_K$  be the  $K$ th roots of unity,

$$t_k = \exp\left(\frac{2\pi i}{K}k\right), \quad k = 1, \dots, K.$$

Define  $g_{K-1}(z)$  as the polynomial

$$g_{K-1}(z) := s_0 z^{K-1} + s_1 z^{K-2} + \dots + s_{K-1}.$$

Note that  $\deg g_{K-1}(z) = K-1$  as  $s_0 \neq 0$ . Without loss of generality we may assume that  $g_{K-1}(t_k) \neq 0$  for  $k = 1, \dots, K$ . (This condition will be needed in Theorem 6.1.)

Let  $w_K(z) := z^K - 1$  and define the symmetric bilinear form  $\langle\langle \cdot, \cdot \rangle\rangle$  as

$$\langle\langle \phi, \psi \rangle\rangle := \sum_{k=1}^K \frac{g_{K-1}(t_k)}{w'_K(t_k)} \phi(t_k) \psi(t_k)$$

for  $\phi, \psi \in \mathcal{P}$ . Note that this form can be evaluated via FFT.

Define the ordinary moments  $\sigma_p$  associated with the form  $\langle\langle \cdot, \cdot \rangle\rangle$  as

$$\sigma_p := \langle\langle 1, z^p \rangle\rangle = \sum_{k=1}^K \frac{g_{K-1}(t_k)}{w'_K(t_k)} t_k^p$$

for  $p = 0, 1, 2, \dots$  and let  $\mathcal{H}_k$  be the  $k \times k$  Hankel matrix

$$\mathcal{H}_k := \left[ \sigma_{p+q} \right]_{p,q=0}^{k-1} = \begin{bmatrix} \sigma_0 & \sigma_1 & \cdots & \sigma_{k-1} \\ \sigma_1 & & \ddots & \vdots \\ \vdots & \ddots & & \vdots \\ \sigma_{k-1} & \cdots & \cdots & \sigma_{2k-2} \end{bmatrix}$$

for  $k = 1, 2, \dots$ . Then the regular FOP  $f_\tau$  of degree  $\tau \geq 1$  associated with the form  $\langle\langle \cdot, \cdot \rangle\rangle$  exists if and only if the matrix  $\mathcal{H}_\tau$  is nonsingular. Also, the following theorem holds (cf. Theorem 2.1).

**THEOREM 6.1.**  $K = \text{rank } \mathcal{H}_{K+p}$  for every nonnegative integer  $p$ .

Thus  $\mathcal{H}_K$  is nonsingular whereas  $\mathcal{H}_\tau$  is singular for  $\tau > K$ . The regular FOP  $f_K$  of degree  $K$  exists while regular FOPs of degree larger than  $K$  do not exist. The polynomial  $f_K$  is easily seen to be

$$f_K(z) = (z - t_1) \cdots (z - t_K) = w_K(z).$$

It is the monic polynomial of degree  $K$  that has  $t_1, \dots, t_K$  as simple zeros.

If  $\mathcal{H}_K$  is strongly nonsingular, then we have a full set  $\{f_0, f_1, \dots, f_K\}$  of regular FOPs. Else, we can proceed in the same way as with the form  $\langle \cdot, \cdot \rangle$ . By filling up the gaps in the sequence of existing regular FOPs it is possible to define a sequence  $\{f_\tau\}_{\tau=0}^\infty$ , with  $f_\tau$  a monic polynomial of degree  $\tau$ , such that if these polynomials are grouped into blocks according to the sequence of regular indices, then polynomials belonging to different blocks are orthogonal with respect to  $\langle \langle \cdot, \cdot \rangle \rangle$ . More precisely, define  $\{f_\tau\}_{\tau=0}^\infty$  as follows. If  $\tau$  is a regular index, then let  $f_\tau$  be the regular FOP of degree  $\tau$ . Else define  $f_\tau$  as  $f_\rho p_{\tau, \rho}$  where  $\rho$  is the largest regular index less than  $\tau$  and  $p_{\tau, \rho}$  is an arbitrary monic polynomial of degree  $\tau - \rho$ . In the latter case  $f_\tau$  is called an *inner polynomial*. If  $p_{\tau, \rho}(z) = z^{\tau - \rho}$  then we say that  $f_\tau$  is defined *by using the standard monomial basis*. The block orthogonality property is expressed by the fact that the Gram matrix  $[\langle \langle f_\tau, f_s \rangle \rangle]_{\tau, s=0}^{K-1}$  is block diagonal. The diagonal blocks are nonsingular, symmetric and zero above the main antidiagonal. If all the inner polynomials in a certain block are defined by using the standard monomial basis, then the corresponding diagonal block has Hankel structure. (Again, see [13] for more details.)

The definition of the form  $\langle \langle \cdot, \cdot \rangle \rangle$  may seem arbitrary. However, there exists a remarkable connection between the forms  $\langle \cdot, \cdot \rangle$  and  $\langle \langle \cdot, \cdot \rangle \rangle$ .

**THEOREM 6.2.** *Let  $\phi, \psi \in \mathcal{P}$ . If  $\deg \phi + \deg \psi \leq K - 1$ , then  $\langle \langle \phi, \psi \rangle \rangle = \langle \phi, \psi \rangle$ .*

**PROOF.** Let  $V(t_1, \dots, t_K)$  be the Vandermonde matrix with nodes  $t_1, \dots, t_K$ ,

$$V(t_1, \dots, t_K) := \begin{bmatrix} 1 & t_1 & \dots & t_1^{K-1} \\ \vdots & \vdots & & \vdots \\ 1 & t_K & \dots & t_K^{K-1} \end{bmatrix}.$$

Then

$$\begin{bmatrix} g_{K-1}(t_1) \\ \vdots \\ g_{K-1}(t_K) \end{bmatrix} = V(t_1, \dots, t_K) \begin{bmatrix} s_{K-1} \\ \vdots \\ s_0 \end{bmatrix}.$$

As  $V(t_1, \dots, t_K)/\sqrt{K}$  is unitary, it follows that

$$\begin{bmatrix} s_{K-1} \\ \vdots \\ s_0 \end{bmatrix} = \frac{1}{K} [V(t_1, \dots, t_K)]^H \begin{bmatrix} g_{K-1}(t_1) \\ \vdots \\ g_{K-1}(t_K) \end{bmatrix}.$$

As  $w_K(z) = z^K - 1$ , it follows that  $w'_K(z) = Kz^{K-1}$  and thus  $w'_K(t_k) = K/t_k$  for  $k = 1, \dots, K$ . Let  $j \in \{1, \dots, K\}$ . Then

$$\sigma_{K-j} = \frac{1}{K} \sum_{k=1}^K g_{K-1}(t_k) t_k^{K-j+1} = \frac{1}{K} \sum_{k=1}^K g_{K-1}(t_k) \overline{t_k^{j-1}}$$

and thus

$$\begin{aligned} \begin{bmatrix} \sigma_{K-1} \\ \sigma_{K-2} \\ \vdots \\ \sigma_0 \end{bmatrix} &= \frac{1}{K} \begin{bmatrix} 1 & \cdots & 1 \\ t_1 & \cdots & t_K \\ \vdots & & \vdots \\ t_1^{K-1} & \cdots & t_K^{K-1} \end{bmatrix} \begin{bmatrix} g_{K-1}(t_1) \\ g_{K-1}(t_2) \\ \vdots \\ g_{K-1}(t_K) \end{bmatrix} \\ &= \frac{1}{K} [V(t_1, \dots, t_K)]^H \begin{bmatrix} g_{K-1}(t_1) \\ \vdots \\ g_{K-1}(t_K) \end{bmatrix} \\ &= \begin{bmatrix} s_{K-1} \\ s_{K-2} \\ \vdots \\ s_0 \end{bmatrix}. \end{aligned}$$

In other words,  $s_p = \sigma_p$  for  $p = 0, 1, \dots, K-1$ . As  $\langle\langle\phi, \psi\rangle\rangle$  depends on  $\sigma_p$  for  $p = 0, 1, \dots, \deg(\phi\psi)$ , this proves the theorem.  $\square$

**COROLLARY 6.3.** *Let  $\tau$  be a nonnegative integer. If  $2\tau - 1 \leq K$ , then  $\tau$  is a regular index for  $\langle\langle\cdot, \cdot\rangle\rangle$  if and only if  $\tau$  is a regular index for  $\langle\cdot, \cdot\rangle$ . Moreover, if  $2\tau \leq K$  and if  $\tau$  is a regular index, then  $f_\tau(z) \equiv \varphi_\tau(z)$ . Else, if  $\tau$  is not a regular index, then  $f_\tau(z) = R_{\tau,\rho}(z)\varphi_\tau(z)$  where  $\rho$  is the largest regular index less than  $\tau$  and  $R_{\tau,\rho}(z)$  is a rational function of type  $[\tau - \rho/\tau - \rho]$ . If  $f_\tau(z)$  and  $\varphi_\tau(z)$  are both defined by using the standard monomial basis, then  $R_{\tau,\rho}(z) \equiv 1$ .*

**COROLLARY 6.4.** *If  $K \geq 2n$  and  $n \leq \tau \leq \lfloor K/2 \rfloor$ , then  $f_\tau(z_k) = 0$  for  $k = 1, \dots, n$  and  $\langle z^p, f_\tau(z) \rangle = 0$  for all  $p \geq 0$ . Also,  $\langle\langle z^p, f_\tau(z) \rangle\rangle = 0$  for  $p = \tau, \dots, K-1-\tau$ . (Note that the latter range may be empty.)*

**COROLLARY 6.5.** *If  $K \geq 2m$  and  $m \leq \tau \leq \lfloor K/2 \rfloor$ , then  $f_\tau(c_j) = \mathcal{O}(\delta^2)$ ,  $\delta \rightarrow 0$  for  $j = 1, \dots, m$  and  $\langle z^p, f_\tau(z) \rangle = \mathcal{O}(\delta^2)$ ,  $\delta \rightarrow 0$  for all  $p \geq \tau$ . Also,  $\langle\langle z^p, f_\tau(z) \rangle\rangle = \mathcal{O}(\delta^2)$ ,  $\delta \rightarrow 0$  for  $p = \tau, \dots, K-1-\tau$ . (Note that the latter range may be empty.)*

Thus, if  $K \geq 2N$ , then we can apply our algorithm of Section 4 to the form  $\langle\langle\cdot, \cdot\rangle\rangle$  and we will obtain exactly the same results as with the form  $\langle\cdot, \cdot\rangle$ . This is an interesting fact in its own right. The main reason, though, that motivated us to introduce the form  $\langle\langle\cdot, \cdot\rangle\rangle$  is the fact that it is related to rational interpolation. We will show that the “denominator polynomials” in a certain linearized rational interpolation problem that is related to the polynomial  $g_{K-1}(z)$  are FOPs with respect to  $\langle\langle\cdot, \cdot\rangle\rangle$ . This will lead to an alternative way to calculate the FOPs  $f_\tau(z)$  and thus, because of Corollary 6.3, the FOPs  $\varphi_\tau(z)$ .

Let  $\sigma$  and  $\tau$  be nonnegative integers such that  $\sigma + \tau + 1 = K$ . Let  $p_\sigma(z)$  and  $q_\tau(z)$  be polynomials, where

$$(6.1) \quad \deg p_\sigma(z) \leq \sigma \quad \text{and} \quad \deg q_\tau(z) \leq \tau,$$



such that the following linearized rational interpolation conditions are satisfied:

$$(6.2) \quad p_\sigma(t_k) - q_\tau(t_k)g_{K-1}(t_k) = 0, \quad k = 1, \dots, K.$$

Each pair of polynomials  $(p_\sigma(z), q_\tau(z))$  that satisfies the degree conditions (6.1) and the interpolation conditions (6.2) is called a *multipoint Padé form* (MPF). The polynomials  $p_\sigma(z)$  and  $q_\tau(z)$  will be called *numerator polynomial* and *denominator polynomial*, respectively.

The interpolation conditions (6.2) lead to a system of  $K$  linear equations in  $K + 1$  unknowns, and thus at least one nontrivial (i.e., whose numerator and denominator polynomial are not identically equal to zero) MPF exists. As (6.2) are homogeneous linear equations, every scalar multiple of a MPF is also a MPF. From now on, we will always assume that MPFs are normalized such that the denominator polynomial is monic. However, the fact that then the number of interpolation conditions is equal to the number of unknown polynomial coefficients, does not guarantee that there exists only one MPF. It merely guarantees that every MPF leads to the same irreducible rational function, called *multipoint Padé approximant* (MPA). Indeed, suppose that there exist two MPAs. Then the numerator polynomial of the difference of these MPAs is a polynomial of degree  $\leq \sigma + \tau$  that vanishes at  $\sigma + \tau + 1$  points. This numerator polynomial is therefore identically equal to zero, which implies that the MPA is unique.

Let  $\mathcal{R}_{\sigma,\tau}$  be the set of rational functions of type  $[\sigma/\tau]$ , i.e., with numerator degree at most  $\sigma$  and denominator degree at most  $\tau$ . A rational interpolation problem that is closely related to (6.2) is the *Cauchy interpolation problem*: find all irreducible rational functions  $r_{\sigma,\tau}(z) \in \mathcal{R}_{\sigma,\tau}$  whose denominator polynomial is monic, such that

$$(6.3) \quad r_{\sigma,\tau}(t_k) = g_{K-1}(t_k), \quad k = 1, \dots, K.$$

This interpolation problem is not always solvable. If a solution exists, then it is unique, and it is equal to the MPA. In general, however, the MPA need not solve the Cauchy interpolation problem: the numerator and denominator polynomials of the MPFs may have common zeros at some interpolation points. The MPA may not satisfy the interpolation condition (6.3) at these points, which are then called *unattainable points*.

Let  $r(z) \in \mathcal{R}_{\sigma,\tau}$  and suppose that  $r(z) = p(z)/q(z)$  where  $p(z)$  and  $q(z)$  are relatively prime polynomials. The *defect* of  $r$  with respect to  $\mathcal{R}_{\sigma,\tau}$  is then defined as

$$\min\{\sigma - \deg p(z), \tau - \deg q(z)\}.$$

The following theorem gives the general solution of the linearized rational interpolation problem.

**THEOREM 6.6.** *The general MPF that corresponds to the degree conditions (6.1) and the interpolation conditions (6.2) is given by*

$$(p_\sigma(z), q_\tau(z)) = (\hat{p}_\sigma(z)s(z)u(z), \hat{q}_\tau(z)s(z)u(z)),$$

where  $\hat{p}_\sigma(z)$ ,  $\hat{q}_\tau(z)$  and  $s(z)$  are uniquely determined polynomials, and where  $u(z)$  is arbitrary. The polynomials  $\hat{p}_\sigma(z)$  and  $\hat{q}_\tau(z)$  are relatively prime, and  $s(z)$  is a divisor of  $w_K(z)$ . Let  $\hat{\delta}_{\sigma,\tau}$  be the defect of  $\hat{p}_\sigma(z)/\hat{q}_\tau(z)$  with respect to  $\mathcal{R}_{\sigma,\tau}$ . Then  $\deg s(z) \leq \hat{\delta}_{\sigma,\tau}$  and  $\deg u(z) \leq \hat{\delta}_{\sigma,\tau} - \deg s(z)$ . The zeros of  $s(z)$  are the unattainable points for the corresponding Cauchy interpolation problem.

PROOF. See, for example, Gutknecht [24, p. 549].  $\square$

The literature on rational interpolation is vast. We will not give a comprehensive account of all the other issues (in particular, the block structure of the Newton–Padé table) that are involved. The reader may wish to consult the papers by Meinguet [41], Antoulas [6, 7, 8], Berrut and Mittelmann [10] or Gutknecht [24, 25, 26, 27], and the references cited therein.

What is of special interest to us, is the fact that the denominator polynomials  $q_\tau(z)$  are formal orthogonal polynomials with respect to  $\langle\langle \cdot, \cdot \rangle\rangle$ .

**THEOREM 6.7.** *Let  $\sigma$  and  $\tau$  be nonnegative integers such that  $\sigma + \tau + 1 = K$ . Let  $(p_\sigma(z), q_\tau(z))$  be a MPF for the degree conditions (6.1) and the interpolation conditions (6.2). Then  $\langle\langle z^p, q_\tau(z) \rangle\rangle = 0$  for  $p = 0, 1, \dots, K - 2 - \deg p_\sigma(z)$  and  $\langle\langle z^p, q_\tau(z) \rangle\rangle \neq 0$  if  $p = K - 1 - \deg p_\sigma(z)$ .*

PROOF. Apparently this orthogonality relation was already known to Jacobi. As it plays a very important role in our paper, we prefer to give a (short, but explicit) proof. See also Egecioğlu and Koç [18] and Gemignani [21] for a slightly weaker version of this theorem.

Let  $p \in \{0, 1, \dots, K - 2 - \deg p_\sigma(z)\}$ . Then

$$(6.4) \quad \sum_{k=1}^K \frac{t_k^p p_\sigma(t_k)}{w'_K(t_k)} = \sum_{k=1}^K \frac{g_{K-1}(t_k)}{w'_K(t_k)} t_k^p q_\tau(t_k)$$

and  $z^p p_\sigma(z)$  is a polynomial of degree  $p + \deg p_\sigma(z) \leq K - 2$ . Lagrange's formula for the polynomial  $y_{K-1}(z)$  of degree  $\leq K - 1$  that interpolates the polynomial  $z^p p_\sigma(z)$  in the points  $t_1, \dots, t_K$  implies that the left-hand side of (6.4) is equal to the coefficient of  $z^{K-1}$  of  $y_{K-1}(z)$ . As  $\deg[z^p p_\sigma(z)] < K - 1$ , it follows that  $y_{K-1}(z) \equiv z^p p_\sigma(z)$  and that the coefficient of  $z^{K-1}$  of  $y_{K-1}(z)$  is equal to zero. It follows that  $\langle\langle z^p, q_\tau(z) \rangle\rangle = 0$  for  $p = 0, 1, \dots, K - 2 - \deg p_\sigma(z)$ . A similar reasoning shows that  $\langle\langle z^p, q_\tau(z) \rangle\rangle \neq 0$  if  $p = K - 1 - \deg p_\sigma(z)$ . This proves the theorem.  $\square$

The following theorem implies that the coefficients (in the standard monomial basis) of the numerator polynomial  $p_\sigma(z)$  of a MPF  $(p_\sigma(z), q_\tau(z))$  that corresponds to the degree conditions (6.1) and the interpolation conditions (6.2) can be expressed as inner products with respect to  $\langle\langle \cdot, \cdot \rangle\rangle$ . This explains how the degree property  $\deg p_\sigma(z) \leq \sigma$  is related to the formal orthogonality property satisfied by  $q_\tau(z)$ .

**THEOREM 6.8.** *Suppose that  $q(z)$  is a polynomial, and let  $p(z)$  be the polynomial of degree  $\leq K - 1$  that interpolates  $g_{K-1}(z)q(z)$  at the points  $t_1, \dots, t_K$ . Let  $p(z) =: p_0 + p_1 z + \dots + p_{K-1} z^{K-1}$ . Then  $p_k = \langle\langle z^{K-1-k}, q(z) \rangle\rangle$  for  $k = 0, 1, \dots, K - 1$ .*

PROOF. The Lagrange representation of  $p(z)$  is given by

$$p(z) = \sum_{k=1}^K \pi_k L_k(z)$$

where

$$\pi_k := \frac{g_{K-1}(t_k)q(t_k)}{w'_K(t_k)} \quad \text{and} \quad L_k(z) := \frac{w_K(z)}{z - t_k}$$

for  $k = 1, \dots, K$ . Let  $L_k(z) =: L_{0,k} + L_{1,k}z + \dots + L_{K-1,k}z^{K-1}$  for  $k = 1, \dots, K$ . Note that  $L_{K-1,1} = \dots = L_{K-1,K} = 1$ . Let

$$V := \begin{bmatrix} 1 & t_1 & \dots & t_1^{K-1} \\ \vdots & \vdots & & \vdots \\ 1 & t_K & \dots & t_K^{K-1} \end{bmatrix}$$

be the Vandermonde matrix with nodes  $t_1, \dots, t_K$ , and let

$$L := \begin{bmatrix} L_{0,1} & \dots & L_{0,K} \\ \vdots & & \vdots \\ L_{K-1,1} & \dots & L_{K-1,K} \end{bmatrix}$$

be the matrix that contains the coefficients of  $L_1(z), \dots, L_K(z)$ . Then

$$\begin{aligned} VL &= \text{diag}(L_1(t_1), \dots, L_K(t_K)) \\ &= \text{diag}(w'_K(t_1), \dots, w'_K(t_K)) \\ &= K \text{diag}(\overline{t_1}, \dots, \overline{t_K}). \end{aligned}$$

As  $V/\sqrt{K}$  is unitary, it follows that  $V^{-1} = V^H/K$ , and thus

$$L = V^H \text{diag}(\overline{t_1}, \dots, \overline{t_K}) = \begin{bmatrix} t_1^{K-1} & \dots & t_K^{K-1} \\ \vdots & & \vdots \\ t_1 & \dots & t_K \\ 1 & \dots & 1 \end{bmatrix}.$$

In other words,  $L_{j,k} = t_k^{K-1-j}$  for  $k = 1, \dots, K$  and  $j = 0, 1, \dots, K-1$ . As  $p_j = \sum_{k=1}^K L_{j,k} \pi_k$  for  $j = 0, 1, \dots, K-1$ , it follows that

$$p_j = \sum_{k=1}^K \frac{g_{K-1}(t_k)}{w'_K(t_k)} t_k^{K-1-j} q(t_k) = \langle\langle z^{K-1-j}, q(z) \rangle\rangle$$

for  $j = 0, 1, \dots, K-1$ . This proves the theorem.  $\square$

The following theorem shows how to construct the sequence of FOPs  $f_\tau(z)$  from the MPFs  $(p_\sigma(z), q_\tau(z))$ . Regular FOPs correspond to denominator polynomials whose degree is equal to  $\tau$  whereas the sizes of the blocks are determined by the actual degrees of the numerator polynomials.

**THEOREM 6.9.** *Let  $\sigma$  and  $\tau$  be nonnegative integers such that  $\sigma + \tau + 1 = K$ . Let  $(p_\sigma(z), q_\tau(z)) = (\hat{p}_\sigma(z)s(z)u(z), \hat{q}_\tau(z)s(z)u(z))$  be the general MPF for the degree conditions (6.1) and the interpolation conditions (6.2), where the polynomials  $\hat{p}_\sigma(z)$ ,  $\hat{q}_\tau(z)$ ,  $s(z)$  and  $u(z)$  are as in Theorem 6.6. If  $\tau$  is a regular index for  $\langle\langle \cdot, \cdot \rangle\rangle$ , then  $\deg(\hat{q}_\tau(z)s(z)) = \tau$ , the FOP  $f_\tau(z)$  is given by  $f_\tau(z) \equiv \hat{q}_\tau(z)s(z)$  and the smallest regular index that is larger than  $\tau$  is equal to  $K - \deg(\hat{p}_\sigma(z)s(z))$ . Conversely, if  $\deg(\hat{q}_\tau(z)s(z)) = \tau$ , then  $\tau$  is a regular index for  $\langle\langle \cdot, \cdot \rangle\rangle$ .*

**PROOF.** Suppose that  $\tau$  is a regular index for  $\langle\langle \cdot, \cdot \rangle\rangle$ . Then  $\det \mathcal{H}_\tau \neq 0$  and there exists precisely one monic polynomial  $f_\tau(z)$  of degree  $\tau$  such that

$$\langle\langle z^p, f_\tau(z) \rangle\rangle = 0 \quad \text{for } p = 0, 1, \dots, \tau - 1.$$

Let  $p(z)$  be the polynomial of degree  $\leq K - 1$  that interpolates  $g_{K-1}(z)f_\tau(z)$  at the points  $t_1, \dots, t_K$ . Then, according to Theorem 6.8,  $\deg p(z) \leq K - \tau - 1 = \sigma$ . Thus  $(p(z), f_\tau(z))$  is a MPF for the degree conditions (6.1) and the interpolation conditions (6.2). In other words, there exists a MPF whose denominator polynomial has degree  $\tau$ . Theorem 6.6 then implies that there exists a monic polynomial  $u_\tau(z)$  of degree  $\tau - \deg(\hat{q}_\tau(z)s(z))$  such that  $f_\tau(z) \equiv \hat{q}_\tau(z)s(z)u_\tau(z)$ . If  $\deg u_\tau(z) > 0$ , then we can choose a different monic polynomial  $\tilde{u}_\tau(z)$  of the same degree. Then  $f_\tau(z) \not\equiv \hat{q}_\tau(z)s(z)\tilde{u}_\tau(z)$  and, by Theorem 6.7,

$$\langle\langle z^p, \hat{q}_\tau(z)s(z)\tilde{u}_\tau(z) \rangle\rangle = 0 \quad \text{for } p = 0, 1, \dots, \tau - 1.$$

As  $\deg(\hat{q}_\tau(z)s(z)\tilde{u}_\tau(z)) = \tau$ , this contradicts the fact that  $f_\tau(z)$  is unique. Thus we may conclude that  $\deg(\hat{q}_\tau(z)s(z)) = \tau$  and  $f_\tau(z) \equiv \hat{q}_\tau(z)s(z)$ . Now Theorem 6.7 implies that

$$\langle\langle z^p, f_\tau(z) \rangle\rangle = 0 \quad \text{for } p = 0, 1, \dots, K - 2 - \deg(\hat{p}_\sigma(z)s(z))$$

and

$$\langle\langle z^p, f_\tau(z) \rangle\rangle \neq 0 \quad \text{if } p = K - 1 - \deg(\hat{p}_\sigma(z)s(z)).$$

The structure of the diagonal blocks of the Gram matrix  $[\langle\langle f_r, f_s \rangle\rangle]_{r,s=0}^{K-1}$  (see, e.g., [13]) then implies that  $\det \mathcal{H}_t = 0$  for  $t = \tau + 1, \dots, K - 1 - \deg(\hat{p}_\sigma(z)s(z))$  and that  $\det \mathcal{H}_t \neq 0$  if  $t = K - \deg(\hat{p}_\sigma(z)s(z))$ .

Suppose that  $\deg(\hat{q}_\tau(z)s(z)) = \tau$ . Then there exists only one MPF for the degree conditions (6.1) and the interpolation conditions (6.2). The polynomial  $\hat{q}_\tau(z)s(z)$  is a monic polynomial of degree  $\tau$  and, according to Theorem 6.7,

$$\langle\langle z^p, \hat{q}_\tau(z)s(z) \rangle\rangle = 0 \quad \text{for } p = 0, 1, \dots, \tau - 1.$$

Suppose that there exists another monic polynomial  $\tilde{f}_\tau(z)$  of degree  $\tau$ ,  $\tilde{f}_\tau(z) \not\equiv \hat{q}_\tau(z)s(z)$ , such that

$$\langle\langle z^p, \tilde{f}_\tau(z) \rangle\rangle = 0 \quad \text{for } p = 0, 1, \dots, \tau - 1.$$

Let  $p(z)$  be the polynomial of degree  $\leq K - 1$  that interpolates  $g_{K-1}(z)\tilde{f}_\tau(z)$  at the points  $t_1, \dots, t_K$ . Then, according to Theorem 6.8,  $\deg p(z) \leq K - \tau - 1 = \sigma$ .

Thus  $(p(z), \tilde{f}_\tau(z))$  is a MPF for the degree conditions (6.1) and the interpolation conditions (6.2). It follows that  $\tilde{f}_\tau(z) \equiv \hat{q}_\tau(z)s(z)$ . In other words, there exists only one monic polynomial of degree  $\tau$  that is orthogonal (with respect to  $\langle\langle \cdot, \cdot \rangle\rangle$ ) to all polynomials of lower degree. Thus  $\tau$  is a regular index for  $\langle\langle \cdot, \cdot \rangle\rangle$ . This proves the theorem.  $\square$

The previous theorem suggests the following look-ahead strategy. Start with  $\tau = 0$  and the corresponding MPF  $(g_{K-1}(z), 1)$ . Then set  $\tau \leftarrow K - \deg g_{K-1}(z)$ . Note that  $\tau = 1$  in case  $\deg g_{K-1}(z) = K - 1$ . Compute the corresponding MPF  $(p_\sigma(z), q_\tau(z))$ . Note that, as  $\tau$  is a regular index for  $\langle\langle \cdot, \cdot \rangle\rangle$ , this MPF is uniquely defined, i.e., the polynomial  $u(z) \equiv 1$  (cf. Theorem 6.6). Use  $K - \deg p_\sigma(z)$  as the next value of  $\tau$ , and so on. Observe that, if  $\deg p_\sigma(z) = \sigma$ , then the next value of  $\tau$  is given by  $\tau + 1$ . The interpolation problems can be solved via the algorithm of Van Barel and Bultheel [51]. This algorithm provides the coefficients of the numerator and the denominator polynomial in the standard monomial basis. It incorporates pivoting.

Of course, in floating-point arithmetic this strategy will only work if one uses a concept of “numerical degree” instead of the classical “degree”. The numerical degree of a polynomial can be defined as follows. Let  $\epsilon > 0$ . The  $\epsilon$ -degree of a polynomial  $p(z) =: p_0 + p_1z + \cdots + p_{K-1}z^{K-1} \in \mathcal{P}$  of degree  $\leq K - 1$  is defined as follows. Let

$$\chi(k) := \frac{\max\{|p_{k+1}|, \dots, |p_{K-1}|\}}{|p_k|}$$

for all  $k \in \{0, 1, \dots, K - 2\}$  such that  $p_k \neq 0$  and  $\chi(k) := \infty$  otherwise. If

$$(6.5) \quad \min_{0 \leq k \leq K-2} \chi(k) \leq \epsilon,$$

then the  $\epsilon$ -degree of  $p(z)$  is defined as the index  $k$  for which the minimum in (6.5) is attained. Else, the  $\epsilon$ -degree of  $p(z)$  is set equal to  $K - 1$ .

The following corollaries provide us with a stopping criterion.

**COROLLARY 6.10.** *Let  $\sigma$  and  $\tau$  be nonnegative integers such that  $\sigma + \tau + 1 = K$ . Let  $(p_\sigma(z), q_\tau(z))$  be a MPF for the degree conditions (6.1) and the interpolation conditions (6.2). If  $K \geq 2n$  and  $n \leq \tau \leq \lfloor K/2 \rfloor$ , then  $\deg p_\sigma(z) \leq \deg q_\tau(z) - 1$ .*

**PROOF.** This follows immediately from Theorem 6.8, Theorem 6.2, and Corollary 6.4.  $\square$

**COROLLARY 6.11.** *Let  $\sigma$  and  $\tau$  be nonnegative integers such that  $\sigma + \tau + 1 = K$ . Let  $(p_\sigma(z), q_\tau(z))$  be a MPF for the degree conditions (6.1) and the interpolation conditions (6.2). Let  $p_\sigma(z) =: p_0 + p_1z + \cdots + p_{K-1}z^{K-1}$ . If  $K \geq 2m$  and  $m \leq \tau \leq \lfloor K/2 \rfloor$ , then*

$$p_k = \mathcal{O}(\delta^2), \quad \delta \rightarrow 0 \quad \text{for} \quad k = \deg q_\tau(z), \dots, K - 1.$$

*In other words, if  $\epsilon$  is sufficiently small, then the  $\epsilon$ -degree of  $p_\sigma(z)$  is less than or equal to  $\deg q_\tau(z) - 1$ .*

**PROOF.** This follows immediately from Theorem 6.8, Theorem 6.2, and Corollary 6.5.  $\square$

In other words, at the end the (numerical) degree of the numerator polynomial is less than or equal to the degree of the denominator polynomial minus one. One can easily verify that this stopping criterion is equivalent to the one used in the algorithm that we have presented in Section 4.

Let us consider the problem of how to evaluate the polynomial  $g_{K-1}(z)$  at the  $K$ th roots of unity  $t_1, \dots, t_K$ . One can easily verify that

$$g_{K-1}(z) = \frac{1}{2\pi i} \int_{\gamma} \frac{t^K - z^K}{t - z} \frac{f'(t)}{f(t)} dt$$

if  $z \notin \gamma$ . Thus, if  $t_k \notin \gamma$  for  $k = 1, \dots, K$ , then

$$g_{K-1}(t_k) = \frac{1}{2\pi i} \int_{\gamma} \frac{t^K - 1}{t - t_k} \frac{f'(t)}{f(t)} dt$$

for  $k = 1, \dots, K$ .

In case  $\gamma$  is the unit circle, one can obtain accurate approximations for

$$g_{K-1}(t_1), \dots, g_{K-1}(t_K)$$

in a very efficient way. Let  $L$  be a positive integer  $\geq K$ , preferably a power of 2. Let  $\omega_1, \dots, \omega_L$  be the  $L$ th roots of unity,

$$\omega_l = \exp\left(\frac{2\pi i}{L}l\right), \quad l = 1, \dots, L.$$

**THEOREM 6.12.** *Suppose that  $\gamma$  is the unit circle. Let  $v_L \in \mathbb{C}^{L \times 1}$  be the vector*

$$v_L := \frac{1}{L} \begin{bmatrix} 1 & \omega_1 & \cdots & \omega_1^{L-1} \\ \vdots & \vdots & & \vdots \\ 1 & \omega_L & \cdots & \omega_L^{L-1} \end{bmatrix}^H \begin{bmatrix} (f'/f)(\omega_1) \\ \vdots \\ (f'/f)(\omega_L) \end{bmatrix}.$$

Then

$$\begin{bmatrix} O_{K \times (L-K)} & I_K \end{bmatrix} v_L \approx \begin{bmatrix} s_{K-1} \\ \vdots \\ s_0 \end{bmatrix}$$

where  $O_{K \times (L-K)}$  denotes the  $K \times (L-K)$  zero matrix and  $I_K$  denotes the  $K \times K$  identity matrix. In other words, the  $K$  last components of  $v_L$  are approximations for  $s_{K-1}, \dots, s_1, s_0$ .

**PROOF.** By approximating

$$s_p = \frac{1}{2\pi i} \int_{\gamma} z^p \frac{f'(z)}{f(z)} dz = \frac{1}{2\pi} \int_0^{2\pi} e^{ip\theta} e^{i\theta} \frac{f'(e^{i\theta})}{f(e^{i\theta})} d\theta, \quad p = 0, 1, 2, \dots,$$

via the trapezoidal rule, we obtain that

$$s_p \approx \frac{1}{L} \sum_{l=1}^L \omega_l^{p+1} \frac{f'(\omega_l)}{f(\omega_l)}, \quad p = 0, 1, 2, \dots$$

It follows that

$$s_p \approx \frac{1}{L} \sum_{l=1}^L \bar{\omega}_l^{L-1-p} \frac{f'(\omega_l)}{f(\omega_l)}, \quad p = 0, 1, \dots, L-1.$$

This proves the theorem.  $\square$

Since

$$\begin{bmatrix} g_{K-1}(t_1) \\ \vdots \\ g_{K-1}(t_K) \end{bmatrix} = \begin{bmatrix} 1 & t_1 & \cdots & t_1^{K-1} \\ \vdots & \vdots & & \vdots \\ 1 & t_K & \cdots & t_K^{K-1} \end{bmatrix} \begin{bmatrix} s_{K-1} \\ \vdots \\ s_0 \end{bmatrix},$$

it follows that we can obtain approximations for  $g_{K-1}(t_1), \dots, g_{K-1}(t_K)$  via one  $L$ -point (inverse) FFT and one  $K$ -point FFT.

## 7 More numerical examples.

We have implemented the strategy described in the previous section in MATLAB. The m-files are available from the authors. In the following examples, the computations have been done via MATLAB 5 (with floating-point relative accuracy  $\approx 2.2204 \cdot 10^{-16}$ ). We have considered the case that  $\gamma$  is the unit circle. Approximations for  $g_{K-1}(t_1), \dots, g_{K-1}(t_K)$  have been computed by using Theorem 6.12. The interpolation problems have been solved via the algorithm of Van Barel and Bultheel [51].

EXAMPLE 7.1. Let us reconsider the problem that we have studied in Example 5.3. As the corresponding  $\gamma$  is given by  $\gamma = \{z \in \mathbb{C} : |z| = 5\}$ , we divide all the zeros by 5 to transform the problem to the unit disk. Recall that  $N = n = 10$  whereas  $m = 4$ . We set  $L = 512$  and  $K = 22$ .

In Figure 7.1 we plot the logarithm with base 10 of the modulus of the coefficients of  $p_\sigma(z)$  for  $\tau = 0, 1, \dots, 5$ . (The logarithm of the modulus of the lowest degree coefficient is shown on the left. In general, the coefficient of  $z^k$  corresponds to the abscis  $k + 1$ .) Note that  $\sigma = 21 - \tau$ .

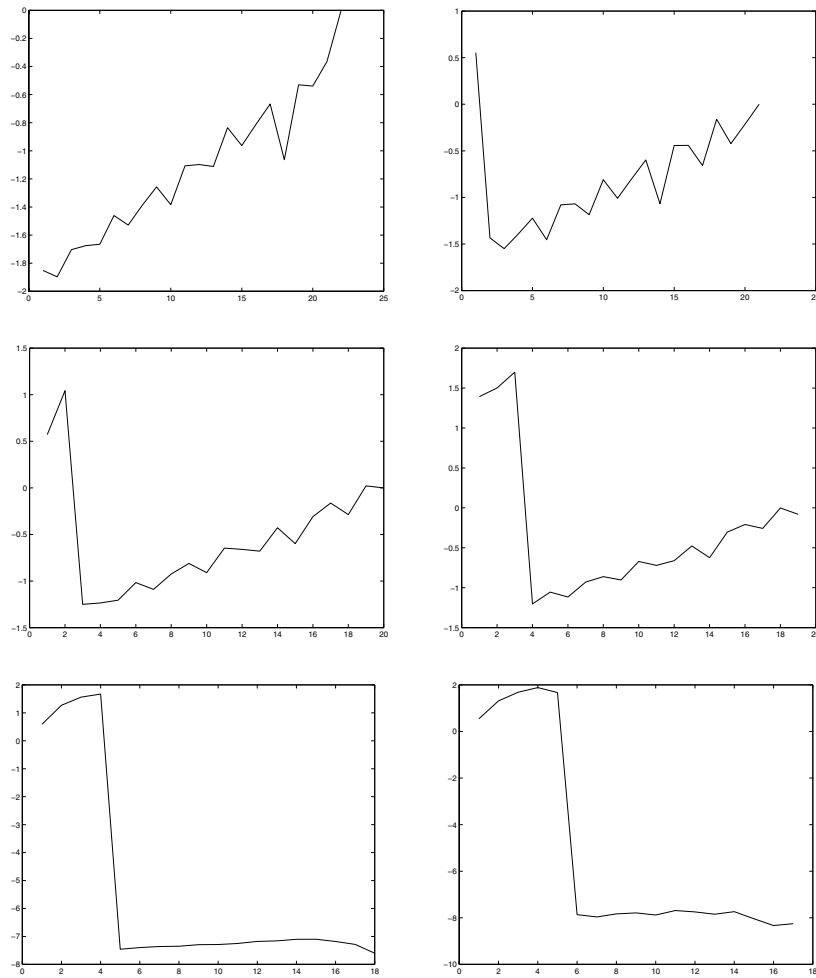
Clearly  $m = 4$ . By multiplying the zeros of  $q_4(z)$ , as computed via the MATLAB command `roots`, by 5 to transform them back to the setting of Example 5.3, we obtain the following:

```
-9.999999564178451 e-01 - i· 5.152536484956870 e-08
 4.000050001653282 e+00 + i· 5.000694740214643 e-05
 2.335838430343232 e-04 + i· 3.000299920075351 e+00
-3.000024926663527 e+00 + i· 3.000149946356137 e+00
```

These values are to be compared with the approximations for the centres of the clusters that we have obtained in Example 5.3, namely

```
-9.999999564181510 e-01 - i· 5.152524762408461 e-08
 4.000050001653271 e+00 + i· 5.000694739720757 e-05
 2.335838430156945 e-04 + i· 3.000299920075392 e+00
-3.000024926663507 e+00 + i· 3.000149946356108 e+00
```

The figure that corresponds to  $\tau = 5$  is included to illustrate that the results given in Corollary 6.11 hold not only for  $\tau = m$  but for  $\tau \geq m$ . The zeros

Figure 7.1: The coefficients of  $p_\sigma(z)$  for  $\sigma = 0, 1, \dots, 5$ .

of  $q_5(z)$  lead to the same approximations for the centres as the zeros of  $q_4(z)$  and one spurious “centre”.

EXAMPLE 7.2. Let

$$f(z) = (\sinh(2z^2) + \sinh(10z) - 1)(\sinh(2z^2) + \sinh(10z) - 1.01) \\ \times (\sinh(2z^2) + \sinh(10z) - 1.02).$$

This function has 21 simple zeros inside the unit circle. They form 7 clusters, where each cluster consists of 3 zeros. Thus  $N = n = 21$  and  $m = 7$ . This example was also studied in [47]. We set  $L = 512$  and  $K = 42$ .



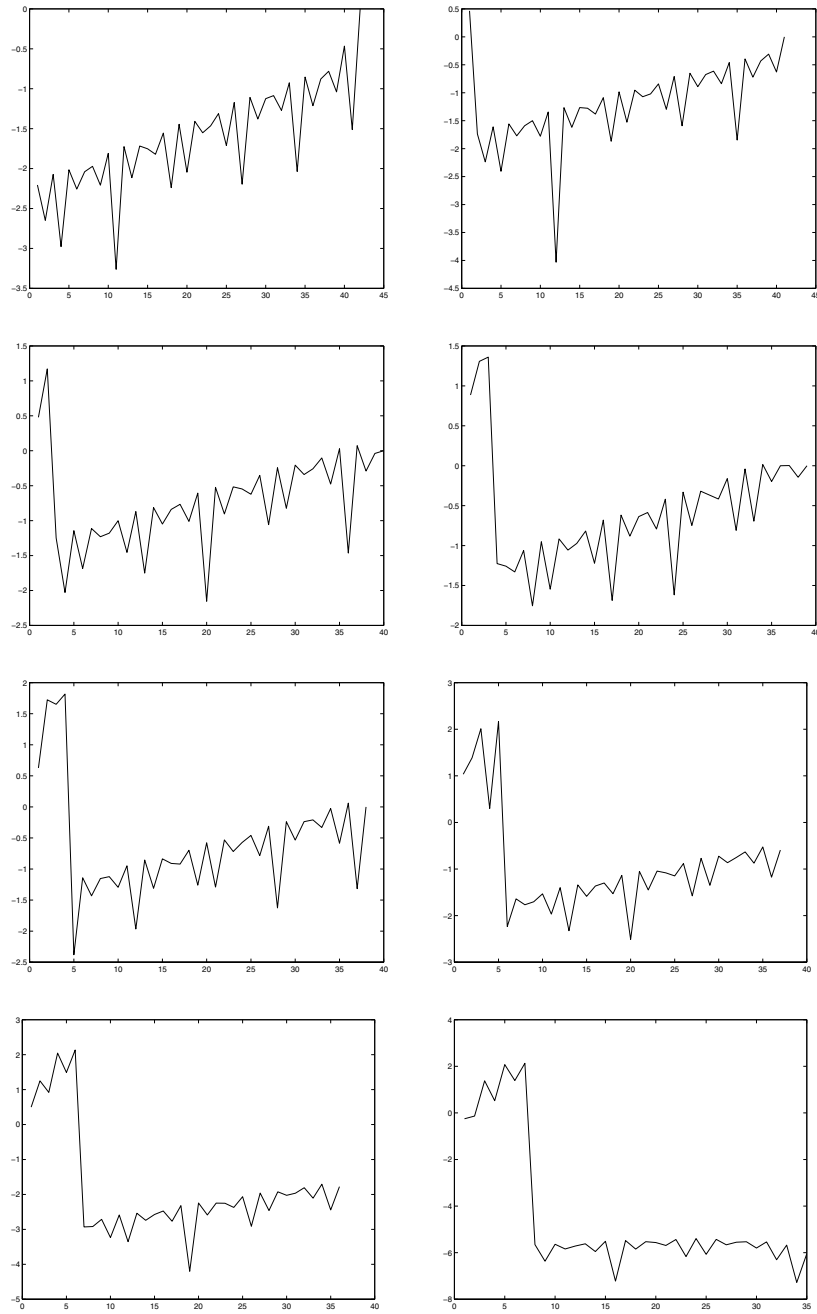


Figure 7.2: The coefficients of  $p_\sigma(z)$  for  $\sigma = 0, 1, \dots, 7$ .

In Figure 7.2 we plot the logarithm with base 10 of the modulus of the coefficients of  $p_\sigma(z)$  for  $\tau = 0, 1, \dots, 7$ .

The zeros of  $q_7(z)$  are given by

$$\begin{aligned} & -1.848537713183581 \text{ e-01} - i \cdot 8.949141853554533 \text{ e-01} \\ & -1.848537713183412 \text{ e-01} + i \cdot 8.949141853554334 \text{ e-01} \\ & -1.003354151041395 \text{ e-01} - i \cdot 3.061151582728444 \text{ e-01} \\ & -1.003354151030711 \text{ e-01} + i \cdot 3.061151582802838 \text{ e-01} \\ & 1.335489810139705 \text{ e-01} - i \cdot 6.084120926164355 \text{ e-01} \\ & 1.335489810131479 \text{ e-01} + i \cdot 6.084120926165633 \text{ e-01} \\ & 8.777826151937687 \text{ e-02} + i \cdot 8.843042856595357 \text{ e-12} \end{aligned}$$

These match the approximations for the centres of the clusters that were given in [47].

### Acknowledgements.

This paper was initiated while Tetsuya Sakurai was staying at the Department of Mathematics and Computer Science of the University of Antwerp. He would like to thank Annie Cuyt for inviting him to Antwerp.

We thank the referees for their very constructive comments and suggestions.

### REFERENCES

1. E. G. Anastasselou, *A formal comparison of the Delves–Lyness and Burniston–Siewert methods for locating the zeros of analytic functions*, IMA J. Numer. Anal., 6 (1986), pp. 337–341.
2. E. G. Anastasselou and N. I. Ioakimidis, *Application of the Cauchy theorem to the location of zeros of sectionally analytic functions*, J. Appl. Math. Phys., 35 (1984), pp. 705–711.
3. E. G. Anastasselou and N. I. Ioakimidis, *A generalization of the Siewert–Burniston method for the determination of zeros of analytic functions*, J. Math. Phys., 25 (1984), pp. 2422–2425.
4. E. G. Anastasselou and N. I. Ioakimidis, *A new method for obtaining exact analytical formulae for the roots of transcendental functions*, Lett. Math. Phys., 8 (1984), pp. 135–143.
5. E. G. Anastasselou and N. I. Ioakimidis, *A new approach to the derivation of exact analytical formulae for the zeros of sectionally analytic functions*, J. Math. Anal. Appl., 112 (1985), pp. 104–109.
6. A. C. Antoulas, *On the scalar rational interpolation problem*, IMA J. Math. Control Inf., 3 (1986), pp. 61–88.
7. A. C. Antoulas, *Rational interpolation and the Euclidean algorithm*, Linear Algebr. Appl., 108 (1988), pp. 157–171.
8. A. C. Antoulas, J. A. Ball, J. Kang, and J. C. Willems, *On the solution of the minimal rational interpolation problem*, Linear Algebr. Appl., 137/138 (1990), pp. 511–573.
9. L. Atanassova, *On the simultaneous determination of the zeros of an analytic function inside a simple smooth closed contour in the complex plane*, J. Comput. Appl. Math., 50 (1994), pp. 99–107.

10. J.-P. Berrut and H. D. Mittelmann, *Matrices for the direct determination of the barycentric weights of rational interpolation*, J. Comput. Appl. Math., 78 (1997), pp. 355–370.
11. A. W. Bojanczyk and G. Heinig, *A multi-step algorithm for Hankel matrices*, J. Complexity, 10 (1994), pp. 142–164.
12. L. C. Botten, M. S. Craig, and R. C. McPhedran, *Complex zeros of analytic functions*, Comput. Phys. Commun., 29 (1983), pp. 245–259.
13. A. Bultheel and M. Van Barel, *Linear Algebra, Rational Approximation and Orthogonal Polynomials*, vol. 6 of Studies in Computational Mathematics, North-Holland, Amsterdam, 1997.
14. E. E. Burniston and C. E. Siewert, *The use of Riemann problems in solving a class of transcendental equations*, Proc. Camb. Philos. Soc., 73 (1973), pp. 111–118.
15. S. Cabay and R. Meleshko, *A weakly stable algorithm for Padé approximants and the inversion of Hankel matrices*, SIAM J. Matrix Anal. Appl., 14 (1993), pp. 735–765.
16. M. P. Carpentier and A. F. D. Santos, *Solution of equations involving analytic functions*, J. Comput. Phys., 45 (1982), pp. 210–220.
17. L. M. Delves and J. N. Lyness, *A numerical method for locating the zeros of an analytic function*, Math. Comp., 21 (1967), pp. 543–560.
18. Ö. Eğecioglu and Ç. K. Koç, *A fast algorithm for rational interpolation via orthogonal polynomials*, Math. Comput., 53 (1989), pp. 249–264.
19. R. W. Freund and H. Zha, *A look-ahead algorithm for the solution of general Hankel systems*, Numer. Math., 64 (1993), pp. 295–321.
20. F. D. Gakhov, *Boundary Value Problems*, vol. 85 of International Series of Monographs in Pure and Applied Mathematics, Pergamon Press, Oxford, 1966.
21. L. Gemignani, *Rational interpolation via orthogonal polynomials*, Comput. Math. Applic., 26 (1993), pp. 27–34.
22. I. Gohberg and I. Koltracht, *Mixed, componentwise, and structured condition numbers*, SIAM J. Matrix Anal. Appl., 14 (1993), pp. 688–704.
23. W. B. Gragg and M. H. Gutknecht, *Stable look-ahead versions of the Euclidean and Chebyshev algorithms*, in Approximation and Computation: A Festschrift in Honor of Walter Gautschi, R. V. M. Zahar, ed., Birkhäuser, Basel, 1994, pp. 231–260.
24. M. H. Gutknecht, *Continued fractions associated with the Newton–Padé table*, Numer. Math., 56 (1989), pp. 547–589.
25. M. H. Gutknecht, *In what sense is the rational interpolation problem well posed?*, Constr. Approx., 6 (1990), pp. 437–450.
26. M. H. Gutknecht, *The rational interpolation problem revisited*, Rocky Mt. J. Math., 21 (1991), pp. 263–280.
27. M. H. Gutknecht, *Block structure and recursiveness in rational interpolation*, in Approximation Theory VII, E. W. Cheney, C. K. Chui, and L. L. Schumaker, eds., Academic Press, New York, 1992, pp. 93–130.
28. V. Hribernic and H. J. Stetter, *Detection and validation of clusters of polynomial zeros*, J. Symbolic Computation, 24 (1997), pp. 667–681.
29. N. I. Ioakimidis, *Quadrature methods for the determination of zeros of transcendental functions—a review*, in Numerical Integration: Recent Developments, Software and Applications, P. Keast and G. Fairweather, eds., Reidel, Dordrecht, 1987, pp. 61–82.

30. N. I. Ioakimidis, *A unified Riemann-Hilbert approach to the analytical determination of zeros of sectionally analytic functions*, J. Math. Anal. Appl., 129 (1988), pp. 134–141.
31. N. I. Ioakimidis, *A note on the closed-form determination of zeros and poles of generalized analytic functions*, Stud. Appl. Math., 81 (1989), pp. 265–269.
32. N. I. Ioakimidis and E. G. Anastasselou, *A new, simple approach to the derivation of exact analytical formulae for the zeros of analytic functions*, Appl. Math. Comp., 17 (1985), pp. 123–127.
33. N. I. Ioakimidis and E. G. Anastasselou, *On the simultaneous determination of zeros of analytic or sectionally analytic functions*, Computing, 36 (1986), pp. 239–247.
34. T. Kailath, *Linear Systems*, Prentice-Hall, Englewood Cliffs, NJ, 1980.
35. P. Kirrinnis, *Newton iteration towards a cluster of polynomial zeros*, in Foundations of Computational Mathematics, F. Cucker and M. Shub, eds., Springer-Verlag, Berlin, 1997, pp. 193–215.
36. P. Kravanja, R. Cools, and A. Haegemans, *Computing zeros of analytic mappings: A logarithmic residue approach*, BIT, 38 (1998), pp. 583–596.
37. P. Kravanja, M. Van Barel, and A. Haegemans, *On computing zeros and poles of meromorphic functions*, in Computational Methods and Function Theory 1977, N. Papamichael, St. Ruscheweyh and E. B. Saff, eds., Series in Approximation and Decompositions, vol. 11, World Scientific, 1999, pp. 359–369.
38. P. Kravanja, M. Van Barel, O. Ragos, M. N. Vrahatis, and F. A. Zafriopoulos, *ZEAL: A mathematical software package for computing zeros of analytic functions*, To appear in Comput. Phys. Commun.
39. T.-Y. Li, *On locating all zeros of an analytic function within a bounded domain by a revised Delves/Lyness method*, SIAM J. Numer. Anal., 20 (1983), pp. 865–871.
40. J. N. Lyness and L. M. Delves, *On numerical contour integration round a closed contour*, Math. Comp., 21 (1967), pp. 561–577.
41. J. Meinguet, *On the solubility of the Cauchy interpolation problem*, in Approximation Theory, A. Talbot, ed., Academic Press, 1970, pp. 137–163.
42. A. Neumaier, *An existence test for root clusters and multiple roots*, Z. Angew. Math. Mech., 68 (1988), pp. 256–257.
43. J. R. Partington, *An Introduction to Hankel Operators*, vol. 13 of London Mathematical Society Student Texts, Cambridge University Press, 1988.
44. M. S. Petković, *Inclusion methods for the zeros of analytic functions*, in Computer Arithmetic and Enclosure Methods, L. Atanassova and J. Herzberger, eds., North-Holland, Amsterdam, 1992, pp. 319–328.
45. M. S. Petković and D. Herceg, *Higher-order iterative methods for approximating zeros of analytic functions*, J. Comput. Appl. Math., 39 (1992), pp. 243–258.
46. M. S. Petković and Z. M. Marjanović, *A class of simultaneous methods for the zeros of analytic functions*, Comput. Math. Appl., 22 (1991), pp. 79–87.
47. T. Sakurai, T. Torii, N. Ohsako, and H. Sugiura, *A method for finding clusters of zeros of analytic function*, in Special Issues of Zeitschrift für Angewandte Mathematik und Mechanik (ZAMM). Issue 1: Numerical Analysis, Scientific Computing, Computer Science, 1996, pp. 515–516. Proceedings of the International Congress on Industrial and Applied Mathematics (ICIAM/GAMM 95), Hamburg, July 3–7, 1995.

- 48. G. W. Stewart, *Perturbation theory for the generalized eigenvalue problem*, in Recent Advances in Numerical Analysis, C. de Boor and G. H. Golub, eds., Academic Press, New York, 1978, pp. 193–206.
- 49. T. Torii and T. Sakurai, *Global method for the poles of analytic function by rational interpolant on the unit circle*, World Sci. Ser. Appl. Anal., 2 (1993), pp. 389–398.
- 50. W. Van Assche, *Orthogonal polynomials in the complex plane and on the real line*, in Special Functions,  $q$ -Series and Related Topics, M. E. H. Ismail, D. R. Masson, and M. Rahman, eds., vol. 14 of Fields Institute Communications, American Mathematical Society, Providence, RI, 1997, pp. 211–245.
- 51. M. Van Barel and A. Bultheel, *A new approach to the rational interpolation problem*, J. Comput. Appl. Math., 32 (1990), pp. 281–289.
- 52. J. H. Wilkinson, *The evaluation of the zeros of ill-conditioned polynomials. Part I*, Numer. Math., 1 (1959), pp. 150–166.
- 53. J.-C. Yakoubsohn, *Approximating the zeros of analytic functions by the exclusion algorithm*, Numerical Algorithms, 6 (1994), pp. 63–88.
- 54. X. Ying and I. N. Katz, *A simple reliable solver for all the roots of a nonlinear function in a given domain*, Computing, 41 (1989), pp. 317–333.