

**NUMERICAL ANALYSIS PROJECT**  
**MANUSCRIPT NA-92-05**

**MAY 1992**

**Adaptive Chebyshev Iterative Methods  
for Nonsymmetric Linear Systems  
Based on Modified Moments**

by

D. Calvetti  
G.H. Golub  
L. Reichel

**NUMERICAL ANALYSIS PROJECT**  
**COMPUTER SCIENCE DEPARTMENT**  
**STANFORD UNIVERSITY**  
**STANFORD, CALIFORNIA 94305**





# Adaptive Chebyshev Iterative Methods for Nonsymmetric Linear Systems Based on Modified Moments

D. Calvetti \*      G.H. Golub †      L. Reichel ‡

## Abstract

Large, sparse nonsymmetric systems of linear equations with a matrix whose eigenvalues lie in the right half plane may be solved by an iterative method based on Chebyshev polynomials for an interval in the complex plane. Knowledge of the convex hull of the spectrum of the matrix is required in order to choose parameters upon which the iteration depends. Adaptive Chebyshev algorithms, in which these parameters are determined by using eigenvalue estimates computed by the power method or modifications thereof, have been described by Manteuffel [16]. This paper presents adaptive Chebyshev iterative methods, in which eigenvalue estimates are computed from modified moments determined during the iterations. The computation of eigenvalue estimates from modified moments requires less computer storage than when eigenvalue estimates are computed by a power method and yields faster convergence for many problems.

## 1 Introduction

The problem of solving a linear system of equations

$$\mathbf{A}\mathbf{x} = \mathbf{b}, \quad \mathbf{A} \in \mathbf{R}^{N \times N}, \quad \mathbf{x}, \mathbf{b} \in \mathbf{R}^N, \quad (1.1)$$

with a large, sparse and nonsymmetric matrix  $\mathbf{A}$  arises in many applications. A Chebyshev iterative method based on scaled Chebyshev polynomials  $p_n$  for an interval in the complex plane can be used to solve (1.1) when the spectrum of  $\mathbf{A}$  lies in the right half plane. This includes matrices with a positive definite symmetric part. Manteuffel [15, 16] discusses such Chebyshev iterative schemes and shows that the iterations depend on two parameters only, the center  $d$  and the focal length  $c$  of an ellipse in the complex plane  $\mathbf{C}$  with foci at  $d \pm c$ . In these schemes, the  $p_n$  are Chebyshev polynomials for the interval between the foci, and are scaled so that  $p_n(\mathbf{O}) = 1$ . The three-term recurrence relation for the  $p_n$  yields an inexpensive recurrence relation for computing a sequence of approximate solutions

---

\*Department of Pure and Applied Mathematics, Stevens Institute of Technology, Hoboken, NJ 07030. Research supported in part by the Design and Manufacturing Institute at Stevens Institute of Technology.

†Department of Computer Science, Stanford University, Stanford, CA 94305-2140. Research supported in part by NSF grant CCR-8821078.

‡Department of Mathematics and Computer Science, Kent State University, Kent, OH 44242. Research supported in part by NSF grant DMS-9002884.

$\mathbf{x}_n$ ,  $n = 1, 2, \dots$ , of (1.1). Let  $\mathbf{x}_0$  denote a given initial approximate solution of (1.1), and introduce the residual vectors  $\mathbf{r}_n := \mathbf{b} - \mathbf{A}\mathbf{x}_n$ ,  $n \geq 0$ . The iterates  $\mathbf{x}_n$  determined by the Chebyshev iterative method are such that

$$\mathbf{e}_n = p_n(\mathbf{A})\mathbf{e}_0, \quad n \geq 0, \quad (1.2)$$

where  $\mathbf{e}_n$  denotes the error in  $\mathbf{x}_n$ , i.e.,

$$\mathbf{e}_n := \mathbf{x}^* - \mathbf{x}_n, \quad \mathbf{x}^* := \mathbf{A}^{-1}\mathbf{b}. \quad (1.3)$$

Let the matrix  $\mathbf{A}$  be diagonalizable and have spectral decomposition

$$\mathbf{A} = \mathbf{W}\mathbf{\Lambda}\mathbf{W}^{-1}, \quad \mathbf{A} = \text{diag}[\lambda_1, \lambda_2, \dots, \lambda_N], \quad \mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N], \quad (1.4)$$

where the eigenvectors  $\mathbf{w}_j$  are scaled so that  $\|\mathbf{w}_j\| = 1$ . Throughout this paper  $\|\cdot\|$  denotes the Euclidean vector norm or the corresponding induced matrix norm. Let  $\mathbf{S}(\mathbf{A})$  denote the spectrum of  $\mathbf{A}$ . It follows from (1.4) that the error  $\mathbf{e}_n$  can be bounded by

$$\|\mathbf{e}_n\| \leq \|\mathbf{W}\| \|p_n(\mathbf{\Lambda})\| \|\mathbf{W}^{-1}\| \|\mathbf{e}_0\| = \|\mathbf{W}\| \|\mathbf{W}^{-1}\| \|\mathbf{e}_0\| \max_{\lambda \in \mathbf{S}(\mathbf{A})} |p_n(\lambda)|. \quad (1.5)$$

We obtain from (1.2) that

$$\mathbf{r}_n = p_n(\mathbf{A})\mathbf{r}_0, \quad n \geq 0, \quad (1.6)$$

and, therefore, a bound similar to (1.5) holds for the residual vectors, also. Because of relation (1.6), the  $p_n$  are sometimes referred to as residual polynomials. If the parameters  $d$  and  $c$  are chosen so that the quantity

$$\max_{\lambda \in \mathbf{S}(\mathbf{A})} |p_n(\lambda)| \quad (1.7)$$

decreases rapidly with  $n$ , then, by (1.5), the norm  $\|\mathbf{e}_n\|$  decreases rapidly as  $n$  increases; see, e.g., [15] for details, where the case when  $\mathbf{A}$  cannot be diagonalized is treated, also. For pronouncedly nonnormal matrices  $\mathbf{A}$ , i.e., when  $\|\mathbf{W}\| \|\mathbf{W}^{-1}\|$  is “huge”, it may be meaningful to consider pseudospectra of  $\mathbf{A}$  instead of the spectrum; see [17] for a discussion. For simplicity, we will in the present paper only discuss convergence in terms of the spectrum  $\mathbf{S}(\mathbf{A})$ . For  $n$  sufficiently large, the scaled Chebyshev polynomials  $p_n$  for the interval between the foci at  $d \pm c$  are of nearly constant magnitude on the boundary of any ellipse, which is not an interval, with foci at  $c \pm d$ . Chebyshev iteration is an attractive solution method if parameters  $d$  and  $c$  exist, such that there is an ellipse with foci at  $d \pm c$  which contains  $\mathbf{S}(\mathbf{A})$  and is not very close to the origin. In particular, this ellipse must not contain the origin. Assuming that such an ellipse exists, its center  $d$  and focal length  $c$  can be determined if  $\mathbf{S}(\mathbf{A})$  is explicitly known. However, in general,  $\mathbf{S}(\mathbf{A})$  is neither known nor easy to determine.

In [16] Manteuffel describes algorithms for dynamic estimation of the parameters  $\mathbf{d}$  and  $c$  based on the power method applied to  $\mathbf{A}$ , or modifications thereof. The parameters  $\mathbf{d}$  and  $c$  are chosen so that  $\mathbf{d} \pm c$  are the foci of the the smallest ellipse containing available estimates of eigenvalues of  $\mathbf{A}$ . As new estimates of the eigenvalues of  $\mathbf{A}$  become available during the iterations, it may be necessary to refit the ellipse so that it encloses all available eigenvalue estimates of  $\mathbf{A}$ . Manteuffel [16] proposes a combinatorial approach for fitting the ellipse. More recently other schemes have also been suggested; see [4, 14]. We review Manteuffel's adaptive Chebyshev algorithms in §2.

A modification of Manteuffel's adaptive schemes is proposed by Elman et al. [7], who replace the power method and its modifications by the Arnoldi process and the GMRES algorithm. The Arnoldi process is applied to compute eigenvalue estimates of  $\mathbf{A}$ , and these estimates are used to determine new parameters  $\mathbf{d}$  and  $c$ . Having computed eigenvalue estimates by the Arnoldi process, the best available approximate solution of (1.1), say  $\mathbf{x}_n$ , can be improved quite inexpensively by the GMRES algorithm before restarting Chebyshev iteration with the new parameters  $\mathbf{d}$  and  $c$ . The scheme proposed by Elman et al. [7] is a hybrid iterative method because it combines Chebyshev iteration with the GMRES algorithm. A recent survey of hybrid iterative schemes can be found in [17].

This paper presents two adaptive Chebyshev algorithms that use modified moments to compute approximations of eigenvalues of  $\mathbf{A}$ . The computed modified moments and the recursion coefficients of the  $\mathbf{p}_n$  are input to the modified Chebyshev algorithm, which determines a nonsymmetric tridiagonal matrix. We compute the eigenvalues of this tridiagonal matrix and consider them as estimates of eigenvalues of  $\mathbf{A}$ . These estimates are used to compute parameters  $\mathbf{d}$  and  $c$  by determining the smallest ellipse that contains the estimates. From the location of the foci at  $\mathbf{d} \pm c$  of this ellipse, the parameters  $\mathbf{d}$  and  $c$  can easily be computed.

The computation of each modified moment requires the evaluation of an inner product of two  $N$ -vectors. The adaptive procedure that we describe in this paper requires  $2\kappa$  modified moments to estimate  $\kappa$  eigenvalues of  $\mathbf{A}$ . The simultaneous calculation of iterates  $\mathbf{x}_n$  and modified moments makes it possible to compute new eigenvalue estimates from modified moments and refit the ellipse in order to determine new values for the parameters  $\mathbf{d}$  and  $c$  as soon as the rate of convergence of the computed residual vectors  $\mathbf{r}_n$  falls below a certain tolerance.

Our numerical experiments indicate that modified moments only have to be computed during the first couple of iterations in order to determine parameters  $\mathbf{d}$  and  $c$  that yield a high rate of convergence. When such parameters have been found, the iterations can proceed without computing further modified moments, and therefore without computing further inner products, until an accurate approximate

solution of (1.1) has been found. Typically, the vast majority of the iterations can be carried out without computing modified moments and ‘inner products. The simplicity of Chebyshev iteration with fixed parameters  $\mathbf{d}$  and  $c$  allows efficient implementations on parallel and vector computers; see Dongarra et al. [5, Chapter 7.1.6] for a recent discussion.

Our schemes for computing modified moments and eigenvalue estimates for a nonsymmetric matrix  $\mathbf{A}$  are extensions of an algorithm described by Golub and Kent [12] for the computation of modified moments and eigenvalue estimates for a symmetric matrix. The computation of modified moments requires the residual vector  $\mathbf{r}_0$  be available. This is the only  $N$ -vector that our adaptive Chebyshev algorithms require stored, in addition to the  $N$ -vectors required by nonadaptive Chebyshev iteration. We note that the adaptive Chebyshev algorithms proposed by Manteuffel and Ashby [16, 2] and by Elman et al. [7] require more  $N$ -vectors to be stored than our schemes. Details of the storage requirements are discussed in §2.

This paper is organized in the following way. In §2 we outline nonadaptive Chebyshev iteration and schemes used by Manteuffel [16] and Elman et al. [7] for determining eigenvalue estimates of  $\mathbf{A}$ . The problem of determining the ellipse that encloses the eigenvalue estimates and yields the smallest convergence factor is treated in §3. This section follows the presentation by Manteuffel [16]. In §4 we discuss how modified moments can be used to reduce the problem of estimating the spectrum of  $\mathbf{A}$  to the computation of the eigenvalues of a certain tridiagonal matrix. This section extends results by Golub and Kent [12]. In §5 we study some properties of modified moments with respect to a complex measure with support in  $\mathbf{C}^2$ , and derive the modified Chebyshev algorithm. Our presentation follows Golub and Gutknecht [12]. We use the modified Chebyshev algorithm to compute the elements of a tridiagonal matrix from the modified moments and the recurrence coefficients of the residual polynomials. These elements are recurrence coefficients of a family of orthogonal polynomials associated with the modified moments. The eigenvalues of this tridiagonal matrix approximate eigenvalues of  $\mathbf{A}$  and are used to determine suitable parameters  $\mathbf{d}$  and  $c$  for Chebyshev iteration.

Manteuffel [16] reports numerical experiments with a modified power method for estimating the spectrum of  $\mathbf{A}$ . In this scheme the power method is applied to a matrix  $\hat{\mathbf{A}}$  obtained by shifting and scaling  $\mathbf{A}$ . Eigenvalue estimates obtained by the power method are known to generally converge most quickly to the eigenvalues of  $\mathbf{A}$  of largest magnitude. The purpose of applying the power method to the shifted and scaled matrix  $\hat{\mathbf{A}}$  is to make eigenvalue estimates important for determining suitable parameters  $\mathbf{d}$  and  $c$  converge quickly to eigenvalues of  $\mathbf{A}$ . In §6 we describe how to use modified moments for  $\hat{\mathbf{A}}$  instead of for  $\mathbf{A}$ . The results of numerical experiments comparing our adaptive schemes

with schemes due to Manteuffel as implemented by Ashby and Manteuffel [2] are presented in §7.

## 2 Adaptive Chebyshev algorithms

In this section we outline the adaptive Chebyshev algorithms by Manteuffel [15, 16] and Elman et al. [7] and introduce notation that will be used in the remainder of the paper. A more detailed discussion of the material presented can be found in [2, 7, 15, 16]. Given the two parameters  $c$  and  $\mathbf{d}$ , Chebyshev iteration for (1.1) can be defined as follows. Let  $\mathbf{x}_0$  be the initial approximate solution, let  $\mathbf{r}_0 := \mathbf{b} - \mathbf{A}\mathbf{x}_0$  and  $\Delta_0 := \frac{1}{\mathbf{d}}\mathbf{r}_0$ . The iterates  $\mathbf{x}_n$  for  $n = 1, 2, \dots$  are defined by

$$\begin{aligned}\mathbf{x}_n &:= \mathbf{x}_{n-1} + \Delta_{n-1}, \\ \mathbf{r}_n &:= \mathbf{b} - \mathbf{A}\mathbf{x}_n, \\ \Delta_n &:= \alpha_n \mathbf{r}_n + \beta_n \Delta_{n-1},\end{aligned}\tag{2.1}$$

where

$$\alpha_n := \frac{2}{c} \frac{T_n\left(\frac{\mathbf{d}}{c}\right)}{T_{n+1}\left(\frac{\mathbf{d}}{c}\right)}, \quad \beta_n := \frac{T_{n-1}\left(\frac{\mathbf{d}}{c}\right)}{T_{n+1}\left(\frac{\mathbf{d}}{c}\right)},\tag{2.2}$$

and  $T_n(\lambda)$  is the Chebyshev polynomial

$$T_n(\lambda) := \cosh\left(n \cosh^{-1}(\lambda)\right).$$

The residual polynomials  $p_n$  in (1.2) and (1.6) are given by

$$p_n(\lambda) = \frac{T_n\left(\frac{\mathbf{d}-\lambda}{c}\right)}{T_n\left(\frac{\mathbf{d}}{c}\right)}.\tag{2.3}$$

Let  $\mathbf{d}$  and  $c$  be the center and focal length, respectively, of the smallest ellipse containing  $\mathbf{S}(\mathbf{A})$ . The assumption that  $\mathbf{S}(\mathbf{A})$  lies in the right half plane and is symmetric with respect to the real axis implies that  $\mathbf{d} > 0$  and  $\mathbf{d}^2 > c^2$ . It therefore suffices to consider pairs of real numbers  $(\mathbf{d}, c^2)$  that lie in

$$\mathcal{R} := \{(\mathbf{d}, c^2) : \mathbf{d} > 0, \mathbf{d}^2 > c^2\}.\tag{2.4}$$

For each  $\lambda \in \mathbb{C}$  define the **asymptotic convergence factor**

$$r(\lambda, \mathbf{d}, c^2) := \left| \frac{\mathbf{d} - \lambda + ((\mathbf{d} - \lambda)^2 - c^2)^{1/2}}{\mathbf{d} + (\mathbf{d}^2 - c^2)^{1/2}} \right|.\tag{2.5}$$

It can be shown that for large  $n$  the component of the error  $\mathbf{e}_n$  in the direction of the eigenvectors  $\mathbf{w}_j$ , cf. (1.4), in each iteration is multiplied by a factor of magnitude roughly equal to  $r(\lambda_j, \mathbf{d}, c^2)$ . Therefore, for sufficiently large  $n$ , the dominating eigenvector components of  $\mathbf{e}_n$  are in the directions of

eigenvectors  $\mathbf{w}_j$  associated with eigenvalues  $\lambda_j$  with largest convergence factor (2.5). From the relation  $\mathbf{r}_n = \mathbf{A}\mathbf{e}_n$ , it follows that for  $n$  sufficiently large the residual error  $\mathbf{r}_n$  also is dominated by the same eigenvector components. It is desirable to choose  $\mathbf{d}$  and  $c$  so that the asymptotic convergence factor (2.5) associated with each eigenvalue of  $\mathbf{A}$  is small. This suggests to let the parameters  $\mathbf{d}$  and  $c$  be the solution of the **mini-max** problem

$$\min_{(d,c^2) \in \mathcal{R}} \max_{\lambda \in \mathcal{S}(\mathbf{A})} r(\lambda, d, c^2). \quad (2.6)$$

The eigenvalues of  $\mathbf{A}$  with largest convergence factors are vertices of the convex hull of  $\mathcal{S}(\mathbf{A})$ . Let  $\mathcal{H}(\mathbf{A})$  denote the set of vertices of the convex hull of  $\mathcal{S}(\mathbf{A})$ . The solution of (2.6) is a function only of the eigenvalues of  $\mathbf{A}$  in  $\mathcal{S}(\mathbf{A})$ . Therefore, Manteuffel's adaptive Chebyshev schemes, as well as our algorithms, seek to determine estimates of the elements of  $\mathcal{H}(\mathbf{A})$ .

In one of the adaptive schemes described in [16], the power method is applied to  $\mathbf{A}$  in order to approximate eigenvalues of the matrix. When the power method is applied to  $\mathbf{A}$ , eigenvalues of large magnitude are typically determined most accurately. Therefore, it may be difficult to determine the vertices of  $\mathcal{H}(\mathbf{A})$  that are closest to the origin with high accuracy in this manner. Manteuffel [16] suggests the following approach to circumvent this problem. Let

$$\mathbf{B}(t) := \frac{\mathbf{d} - z + ((\mathbf{d} - z)^2 - c^2)^{1/2}}{\mathbf{d} + (d^2 - c^2)^{1/2}} \quad (2.7)$$

and apply the power method to  $\mathbf{B}(\mathbf{A})$  in order to estimate the eigenvalues of largest magnitude of this operator. If  $b_j$  is an eigenvalue of  $\mathbf{B}(\mathbf{A})$  and

$$g := \mathbf{d} + (d^2 - c^2)^{1/2}, \quad (2.8)$$

then

$$\lambda_j := d - \frac{1}{2} \left( gb_j + \frac{c^2}{gb_j} \right) \quad (2.9)$$

is an eigenvalue of  $\mathbf{A}$  and  $r(\lambda_j, d, c^2) = |B(b_j)|$ . Therefore, the power method applied to  $\mathbf{B}(\mathbf{A})$  typically yields highest accuracy for eigenvalues of  $\mathbf{A}$  with largest convergence factor. In order to avoid the nonlinearity of the relation (2.9) between  $\lambda_j$  and  $b_j$ , Manteuffel [16] also proposed to apply the power method to the matrix

$$\hat{\mathbf{A}} := 2g(d\mathbf{I} - \mathbf{A}) = g^2\mathbf{S} - c^2\mathbf{S}^{-1}. \quad (2.10)$$

The linearity of the relation between  $\mathbf{A}$  and  $\hat{\mathbf{A}}$  makes it simple to compute the eigenvalues of  $\mathbf{A}$  from those of  $\hat{\mathbf{A}}$ , and, moreover, the power method applied to  $\hat{\mathbf{A}}$  typically yields highest accuracy for eigenvalues with largest convergence factor.



In the implementation of adaptive Chebyshev iteration by Ashby and Manteuffel [2], the power method is applied to  $\mathbf{A}$ ,  $\mathbf{B}(\mathbf{A})$  or  $\hat{\mathbf{A}}$  in order to determine four estimates of eigenvalues. The implementation [2] based on the power method applied to  $\mathbf{A}$  or  $\mathbf{B}(\mathbf{A})$  requires the storage of four N-vectors in addition to the vectors used for nonadaptive Chebyshev iteration (2.1). The implementation [2] based the power method applied to  $\hat{\mathbf{A}}$  requires the storage of five additional N-vectors. The adaptive schemes based on the different power methods also differ in their operation count. The power method applied to  $\mathbf{A}$  is, generally, the most expensive scheme; it requires four matrix-vector products with the matrix  $\mathbf{A}$  for the computation of each set of four eigenvalue estimates, in addition to the matrix-vector products needed to compute the iterates  $\mathbf{x}_n$  by Chebyshev iteration.

An alternative to the power methods for computing estimates of eigenvalues of  $\mathbf{A}$  is provided by the Arnoldi process [1]. This method is a Galerkin scheme for approximating  $m$  eigenvalues of  $\mathbf{A}$ , in which the test and trial spaces are a Krylov subspace  $K_m(\mathbf{A}, \mathbf{v}) := \text{span}\{\mathbf{v}, \mathbf{A}\mathbf{v}, \dots, \mathbf{A}^{m-1}\mathbf{v}\}$ , where  $\mathbf{v}$  is a vector in  $\mathbf{R}^N$ . The scheme requires the storage of an orthonormal basis of  $K_m(\mathbf{A}, \mathbf{v})$ , and the operation count is  $\mathcal{O}(m^2N)$ . In the application of the Arnoldi process to computing eigenvalue estimates for Chebyshev iteration discussed in [7]  $m$  is chosen to be four. See [7] for further details.

### 3 Fitting of the ellipse

In this section we discuss how to compute the parameters  $\mathbf{d}$  and  $c$  for Chebyshev iteration. Let  $\mathcal{F}(d, c, \mathbf{a})$  denote the ellipse with center  $\mathbf{d}$ , focal length  $c$  and semi-major axis  $\mathbf{a} \geq \mathbf{0}$ , i.e.,

$$\mathcal{F}(d, c, a) := \{z \in \mathbf{C} : |z - d + c| + |z - d - c| \leq 2a\}.$$

Let  $\Pi_n$  be the set of polynomials of degree at most  $n$ , and also define the subset  $\tilde{\Pi}_n := \{q : q \in \Pi_n, q(0) = 1\}$ . The following theorem shows that the residual polynomials  $p_n$  for Chebyshev iteration with parameters  $\mathbf{d}$  and  $c$  minimize the limit as  $n \rightarrow \infty$  of the  $n^{\text{th}}$  root of the quantity (1.7) when  $\mathcal{S}(\mathbf{A}) := \mathcal{F}(d, c, \mathbf{a})$ .

**Theorem 1** ([15]) *Assume that  $\mathbf{0} \notin \mathcal{F}(d, c, \mathbf{a})$ . Let  $t_n \in \tilde{\Pi}_n$  satisfy*

$$\max_{\lambda \in \mathcal{F}(d, c, \mathbf{a})} |t_n(\lambda)| = \min_{q \in \tilde{\Pi}_n} \max_{\lambda \in \mathcal{F}(d, c, \mathbf{a})} |q(\lambda)| ,$$

*and let  $p_n$  be given by (2.3). Then*

$$\lim_{n \rightarrow \infty} \left[ \max_{\lambda \in \mathcal{F}(d, c, \mathbf{a})} |t_n(\lambda)| \right]^{1/n} = \lim_{n \rightarrow \infty} \left[ \max_{\lambda \in \mathcal{F}(d, c, \mathbf{a})} |p_n(\lambda)| \right]^{1/n} . \quad (3.1)$$

*Moreover,*

$$\lim_{n \rightarrow \infty} |p_n(\lambda)|^{1/n} = r(\lambda, d, c^2)$$

for any  $\lambda \in \mathcal{C}$ , and, in particular,

$$\lim_{n \rightarrow \infty} |p_n(\lambda)|^{1/n} = r(d, c, a, \mathbf{d}, c^2) < 1 \quad (3.2)$$

for any  $\lambda$  on the boundary of  $\mathcal{F}(d, c, a)$ . □

In the terminology used, e.g., in [15, 20], equality (3.1) shows that the residual polynomials  $p_n$  given by (2.3) yield an asymptotically optimal rate of convergence with respect to  $\mathcal{F}(d, c, a)$ . We note that this result can be improved in several ways. Since the polynomials  $p_n$  are scaled Faber polynomials for  $\mathcal{F}(d, c, a)$  results by Eiermann [6] on Faber polynomials imply that

$$\max_{\lambda \in \mathcal{F}(d, c, a)} |p_n(\lambda)| / \max_{\lambda \in \mathcal{F}(d, c, a)} |t_n(\lambda)| \leq \gamma \log n, \quad n \geq 1, \quad (3.3)$$

for some constant  $\gamma$  independent of  $n$ . Moreover, Fischer and Freund [8] have recently shown that for many ellipses the right hand side of (3.3) can be replaced by one.

Formula (3.2) shows that, if  $\mathbf{d}$  and  $c$  are the center and focal length of a small ellipse enclosing the spectrum of  $\mathbf{A}$ , and if this ellipse is not very close to the origin, then  $\max_{\lambda \in \mathcal{S}(\mathbf{A})} |p_n(\lambda)|$  converges rapidly to zero as  $n$  increases. Moreover, if the matrix  $\mathbf{A}$  is not very far from normal, then (1.5) shows that the norm  $\|\mathbf{e}_n\|$  also converges rapidly to zero with increasing  $n$ .

We now outline the scheme of Manteuffel [15] for computing the best ellipse with respect to the spectrum  $\mathcal{S}(\mathbf{A})$ , i.e., we compute the best parameters  $\mathbf{d}$  and  $c$  for Chebyshev iteration, when the spectrum  $\mathcal{S}(\mathbf{A})$  is given. However, we remark that when carrying out Chebyshev iteration,  $\mathcal{S}(\mathbf{A})$  is generally not known. The adaptive Chebyshev algorithm therefore computes the best ellipse with respect to a set of eigenvalue estimates computed during the iterations. This set is typically updated a few times during the iterations.

Since  $\mathbf{A}$  is real, the set  $\mathcal{H}(\mathbf{A})$  is symmetric with respect to the real axis, and the foci of the smallest ellipse enclosing  $\mathcal{H}(\mathbf{A})$  are either real or complex conjugate. The center  $\mathbf{d}$  and focal length  $c$  of the smallest ellipse containing  $\mathcal{H}(\mathbf{A})$  are such that  $(d, c^2) \in \mathcal{R}$ . Thus, the mini-max problem (2.6) can be replaced by the simpler mini-max problem

$$\min_{(d, c^2) \in \mathcal{R}} \max_{\lambda \in \mathcal{H}^+(\mathbf{A})} r(\lambda, \mathbf{d}, c^2), \quad (3.4)$$

where  $\mathcal{H}^+(\mathbf{A}) := \{\lambda \in \mathcal{H}(\mathbf{A}) : \text{Im}(\lambda) \geq 0\}$ . The theorem below is helpful for the solution of (3.4).

**Theorem 2** ([3, 15]) *Let the set  $M \subset \mathcal{R}^2$  be closed and bounded, and let  $\mathcal{S}(\mathbf{A}) = \{\lambda_i\}_{i=1}^N$ . Then  $\{r(\lambda_i, d, c^2)\}_{i=1}^N$  is a finite set of real-valued functions of two variables  $(d, c^2)$ , continuous on  $M$ . Let*

$m(d, c^2) := \max_i r(\lambda_i, d, c^2)$ . Then  $m(d, c^2)$  has a minimum at some point  $(d_0, c_0^2)$ . If  $(d_0, c_0^2)$  is in the interior of  $M$  then one of the following statements holds:

1. The point  $(d_0, c_0^2)$  is a local minimum of  $r(\lambda_i, d, c^2)$  for some  $i$  such that  $r(\lambda_i, d_0, c_0^2) = m(d_0, c_0^2)$ .
2. The point  $(d_0, c_0^2)$  is a local minimum among the loci  $\{(d, c^2) \in M : r(\lambda_i, d, c^2) = r(\lambda_j, d, c^2)\}$  for some  $i$  and  $j$  such that  $m(d_0, c_0^2) = r(\lambda_j, d_0, c_0^2) = r(\lambda_i, d_0, c_0^2)$ .
3. The point  $(d_0, c_0^2)$  is such that for some  $i, j$  and  $k$ ,  $m(d_0, c_0^2) = r(\lambda_i, d_0, c_0^2) = r(\lambda_j, d_0, c_0^2) = r(\lambda_k, d_0, c_0^2)$ .

□

Manteuffel [15] presents an algorithm for the solution of (3.4) based on the following observations that are a consequence of Theorem 2.

1. If  $\mathcal{H}^+(A) = \{\lambda_1\}$ , then  $\mathbf{d} = \mathbf{x}_1$  and  $c^2 = -y_1^2$ , where  $\lambda_1 = \mathbf{x}_1 + iy_1$ ,  $i := \sqrt{-1}$ .
2. If  $\mathcal{H}^+(A) = \{\lambda_1, \lambda_2\}$ , then the optimal parameters  $(\mathbf{d}, c^2)$  correspond to a point lying on the intersection of the two surfaces

$$r(\lambda_1, \mathbf{d}, c^2) = r(\lambda_2, \mathbf{d}, c^2).$$

The point corresponding to the optimal parameters when the positive convex hull contains only two points is called the *pairwise best point*, and the associated ellipse passing through  $\lambda_1$  and  $\lambda_2$  is called the *pairwise best ellipse*.

3. If  $\mathcal{H}^+(A)$  contains three or more eigenvalues, then the solution to (3.4) must be either a pairwise best point or it is the intersection of three surfaces. Given the pairwise best point for two eigenvalues  $\lambda_1$  and  $\lambda_2$ , this is the best point if the associated pairwise best ellipse contains all eigenvalues in the closure of its interior. If no pairwise best point is the solution to (3.4), then determine the *three-way point* on the intersection of the three surfaces:

$$r(\lambda_1, \mathbf{d}, c^2) = r(\lambda_2, \mathbf{d}, c^2) = r(\lambda_3, \mathbf{d}, c^2)$$

and the associated *three-way ellipse*. If the associated three-way ellipse contains all eigenvalues of  $A$  in the closure of its interior then the three-way point is a feasible point. The three-way feasible point with smallest convergence factor is the solution to the mini-max problem (3.4).

A detailed description, and some further simplifications, of the scheme for fitting the ellipse outlined above are presented in [15].

## 4 Modified moments

In this section we define moments and modified moments and discuss how they can be used to gain spectral information about the matrix  $\mathbf{A}$  while computing approximate solutions  $\mathbf{x}_n$  of (1.1) by Chebyshev iteration (2.1). Let  $\mathbf{A}$  have spectral resolution (1.4), let  $p_n$  be residual polynomials (2.3) and express  $\mathbf{r}_0$  in the basis of eigenvectors  $\{\mathbf{w}_j\}_{j=1}^N$ , i.e.,

$$\mathbf{r}_0 = \sum_{i=1}^N \alpha_i \mathbf{w}_i.$$

Then it follows from  $\mathbf{p}(\mathbf{A}) = W p_n(\Lambda) W^{-1}$  that

$$\mathbf{r}_n = \sum_{i=1}^N \alpha_i p_n(\lambda_i) \mathbf{w}_i.$$

Introduce the inner product

$$(\mathbf{r}_k, \mathbf{r}_l) := \mathbf{r}_k^T \mathbf{r}_l = \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \mathbf{w}_i^T \mathbf{w}_j p_k(\lambda_i) p_l(\lambda_j),$$

and let  $\gamma(\lambda, \eta)$  be the complex symmetric measure with support in  $\mathbf{C}^2$  and with ‘jumps’ of height  $\gamma_{ij} := \alpha_i \alpha_j \mathbf{w}_i^T \mathbf{w}_j$  at the points  $(\lambda_i, \lambda_j) \in \mathbf{C}^2$ , for  $\lambda_i, \lambda_j \in \mathbf{S}(\mathbf{A})$ . Then the inner product can be written as

$$(\mathbf{r}_k, \mathbf{r}_l) = \sum_{i=1}^N \sum_{j=1}^N p_k(\lambda_i) p_l(\lambda_j) \gamma_{ij} = \int_{\mathbf{C}} \int_{\mathbf{C}} p_k(\lambda) p_l(\eta) d\gamma(\lambda, \eta). \quad (4.1)$$

We remark that if the matrix  $\mathbf{A}$  is symmetric, then its eigenvalues are all real and we can choose the eigenvector matrix  $W$  to be orthogonal. In this case (4.1) simplifies to

$$(\mathbf{r}_k, \mathbf{r}_l) = \sum_{i=1}^N \alpha_i^2 p_k(\lambda_i) p_l(\lambda_i) = \int_{\mathbf{R}} p_k(\lambda) p_l(\lambda) d\alpha(\lambda), \quad (4.2)$$

where  $\alpha$  is a real measure with support on  $\mathbf{R}$  and with jumps of height  $\alpha_i^2$  at  $\lambda_i \in \mathbf{S}(\mathbf{A})$ .

Let

$$\gamma_i := \sum_{j=1}^N \gamma_{ij}, \quad 1 \leq i \leq N,$$

and introduce the linear functional  $\phi$  associated with the measure  $\gamma$  by

$$\phi(q) := \sum_{i=1}^N \tilde{\gamma}_i q(\lambda_i), \quad (4.3)$$

where  $q$  is a polynomial. We are now in a position to define **moments** associated with the measure  $\gamma$  by

$$\mu_k := \phi(\lambda^k), \quad k = 0, 1, 2, \dots, \quad (4.4)$$

as well as **modified moments** with respect to the residual polynomials

$$\nu_k := \phi(p_k), \quad \mathbf{k} = 0, 1, 2, \dots \quad (4.5)$$

For future reference, we note that

$$\nu_k = (\mathbf{r}_k, \mathbf{r}_0), \quad k = 0, 1, 2, \dots \quad (4.6)$$

Let  $\ll \cdot, \cdot \gg$  denote the bilinear form generated by  $\phi$ ,

$$\ll f, g \gg := \phi(fg), \quad (4.7)$$

where  $f$  and  $g$  are polynomials, and assume that there is a family of **monic** formal orthogonal polynomials  $\{\pi_i\}_{i=0}^N$  associated with the measure  $\gamma$ , i.e.,

$$\ll \pi_k, \pi_l \gg \begin{cases} = 0, & \text{if } k \neq l, 0 \leq k, l \leq N, \\ \neq 0, & \text{if } k = l, 0 \leq k < N. \end{cases}$$

The bilinear form (4.7) has the property that

$$\ll zf, g \gg = \ll f, zg \gg,$$

and this implies that the  $\pi_k$  satisfy a three-term recurrence relation

$$\begin{aligned} \pi_{k+1}(z) &= (z - \hat{\alpha}_k)\pi_k(z) - \hat{\beta}_k\pi_{k-1}(z), & 0 \leq k < N, \\ \pi_0(z) &= 1, & \pi_{-1}(z) = 0, \end{aligned} \quad (4.8)$$

see, e.g., [9] for a proof. It follows from (4.3) that the zeros of  $\pi_N(z)$  are the eigenvalues of the matrix  $\mathbf{A}$ . The eigenvalues of  $\mathbf{A}$  can therefore be computed as the eigenvalues of the tridiagonal matrix defined by the coefficients in the three-term recurrence relation for the  $\pi_k$ . In the following two sections we show how to construct the tridiagonal matrix containing the recurrence coefficients for the polynomials  $\pi_k$  from the modified moments and the parameters  $\mathbf{d}$  and  $\mathbf{c}$  of Chebyshev iteration. A discussion based on formula (4.2) on how modified moments can be applied to estimate eigenvalues of a symmetric positive definite matrix can be found in [12].

## 5 Computing eigenvalue estimates

This section derives a matrix identity that connects modified moments associated with a complex measure with the recurrence coefficients of a family of orthogonal polynomials associated with this measure. This identity is the basis of the modified Chebyshev algorithm for computing recursion

coefficients for the orthogonal polynomials, and has been discussed in more detail, also including degenerate cases, by Golub and Gutknecht [11]. The recursion coefficients determine a **tridiagonal** matrix, whose eigenvalues are estimates of eigenvalues of  $\mathbf{A}$ . We discuss the computation of these estimates in the end of this section. Throughout this section we assume that  $N = \infty$  in the formulas of §4, and, in particular, that the measure  $\gamma$  has infinitely many points of support, that the moments  $\mu_k$  given by (4.4) are defined and finite for all  $k \geq 0$ , and that there is a complete family of **monic** formal orthogonal polynomials  $\{\pi_k\}_{k=0}^{\infty}$  associated with  $\gamma$ .

Let  $\{\tau_n\}_{n=0}^{\infty}$  be a family of polynomials that satisfy a three-term recurrence relation. In particular, we are interested in the case when the  $\tau_n$  are the residual polynomials  $p_n$ , given by (2.3), for Chebyshev iteration. It is easy to show that the  $p_n$  satisfy the three-term recurrence relation

$$\begin{aligned} \frac{1}{\alpha_n} p_{n+1}(\lambda) &= \left( \frac{\beta_n}{\alpha_n} + \frac{1}{\alpha_n} - \lambda \right) p_n(\lambda) - \frac{\beta_n}{\alpha_n} p_{n-1}(\lambda), \quad n \geq 1, \\ p_1(\lambda) &= 1 - \frac{\lambda}{d}, \quad p_0(\lambda) = 1, \end{aligned} \quad (5.1)$$

where the coefficients  $\alpha_n$  and  $\beta_n$  are defined by (2.2). Introduce the quantities

$$\sigma_{mn} := \phi(\tau_m \pi_n).$$

It follows from (4.5) that, if  $\tau_k = p_k$ , then

$$\sigma_{m0} = \nu_m, \quad m = 0, 1, 2, \dots, \quad (5.2)$$

and the orthogonality of the  $\pi_n$  yields  $\phi(\tau_m \pi_n) = 0$  for  $m < n$ , thus

$$\sigma_{mn} = 0 \text{ for } m < n. \quad (5.3)$$

In order to derive matrix relations that will be used in the calculation of the recurrence coefficients  $\hat{\alpha}_k$  and  $\hat{\beta}_k$  for the polynomials  $\pi_k$ , we introduce the following semi-infinite vectors

$$\boldsymbol{\pi} = [\pi_0, \pi_1, \dots]^T, \quad \boldsymbol{\tau} = [\tau_0, \tau_1, \dots]^T,$$

and semi-infinite matrices

$$\mathbf{H} := \begin{bmatrix} \hat{\alpha}_0 & \beta_1 & & & 0 \\ 1 & \hat{\alpha}_1 & \beta_2 & & \\ & 1 & \hat{\alpha}_2 & \hat{\beta}_3 & \\ 0 & & \ddots & \ddots & \ddots \end{bmatrix}, \quad \mathbf{T} := \begin{bmatrix} \tau_{00} & \tau_{01} & & & 0 \\ \tau_{10} & \tau_{11} & \tau_{12} & & \\ & \tau_{21} & \tau_{22} & \ddots & \\ 0 & & \ddots & \ddots & \end{bmatrix},$$

$$\mathbf{S} := \begin{bmatrix} \sigma_{00} & & & & 0 \\ \sigma_{10} & \sigma_{11} & & & \\ \sigma_{20} & \sigma_{21} & \sigma_{22} & & \\ \vdots & \vdots & \vdots & \ddots & \end{bmatrix},$$

where  $\mathbf{H}$  and  $\mathbf{T}$  are tridiagonal and  $\mathbf{S}$  is lower triangular. The nonvanishing entries of the matrix  $\mathbf{H}$  are recurrence coefficients for the polynomials  $\pi_k$ , see (4.8), and are to be computed. The nontrivial entries  $\tau_{jk}$  of the matrix  $\mathbf{T}$  are recurrence coefficients of the polynomials  $\tau_j$  and are assumed to be explicitly known. In particular, if  $\tau_j = p_j$ , then we obtain from (5.1) that

$$\begin{aligned} \tau_{00} &= d, & \tau_{10} &= -d, \\ \tau_{n-1,n} &= -\frac{p_n}{\alpha_n}, & \tau_{nn} &= \frac{p_n + 1}{\alpha_n}, & \tau_{n+1,n} &= -\frac{1}{\alpha_n}, & n &\geq 1. \end{aligned}$$

We write the three-term recurrence relations for the  $\pi_k$  and  $\tau_k$  in the form

$$z\pi^T(z) = \pi(z)^T H, \quad z\tau^T(z) = \tau^T(z) T. \quad (5.4)$$

Define the functional  $\hat{\phi}$  on the set of vectors of polynomials by

$$\hat{\phi}([q_0, q_1, \dots]^T) := [\phi(q_0), \phi(q_1), \dots]^T, \quad q_n \in \Pi_n.$$

Applying  $\hat{\phi}$  to the rank-one matrix  $\tau\pi^T$  yields  $\hat{\phi}(\tau\pi^T) = \mathbf{S}$ , and it follows from (5.4) that

$$\mathbf{S}\mathbf{H} = \hat{\phi}(\tau\pi^T)\mathbf{H} = \hat{\phi}(\tau\pi^T\mathbf{H}) = \hat{\phi}(\tau z\pi^T) = \hat{\phi}((z\tau^T)^T\pi^T) = \hat{\phi}(T^T\tau\pi^T) = T^T\hat{\phi}(\tau\pi^T) = T^T\mathbf{S}. \quad (5.5)$$

This matrix identity is the basis of the **modified Chebyshev algorithm**, described in [10, 11, 18, 21], for computing the recurrence coefficients for the polynomials  $\pi_k$  from the recurrence coefficients for the  $\tau_k$  and the modified moments  $\nu_k$ . Let  $\mathbf{H}_\kappa$  denote the  $\kappa \times \kappa$  leading principal submatrix of  $\mathbf{H}$ . We derive the modified Chebyshev algorithm for computing the entries of  $\mathbf{H}_\kappa$ . Equating elements in the left and right hand sides of equation (5.5) yields

$$\sigma_{i,j+1} + \hat{\alpha}_j\sigma_{ij} + \hat{\beta}_j\sigma_{i,j-1} = \tau_{i-1,i}\sigma_{i-1,j} + \tau_{ii}\sigma_{ij} + \tau_{i+1,i}\sigma_{i+1,j}. \quad (5.6)$$

If  $i < j - 1$ , then both the right hand side and left hand side of (5.6) vanish, because  $\mathbf{S}$  is lower triangular. When  $i = j - 1$ , formula (5.6) yields

$$\hat{\beta}_j\sigma_{j-1,j-1} = \tau_{j,j-1}\sigma_{jj}, \quad (5.7)$$

and for  $i = j$  we obtain

$$\hat{\alpha}_j\sigma_{jj} + \hat{\beta}_j\sigma_{j,j-1} = \tau_{jj}\sigma_{jj} + \tau_{j+1,j}\sigma_{j+1,j}. \quad (5.8)$$

The coefficients  $\hat{\alpha}_j$  and  $\hat{\beta}_j$  are computed by (5.7) and (5.8), and we note that this requires only the diagonal and subdiagonal elements of the matrix  $\mathbf{S}$ . These entries of  $\mathbf{S}$  can be generated recursively

starting from  $\{\sigma_{j0}\}_{j=0}^{2\kappa-1}$  defined by (5.2). The computations proceed as follows. Initialize

$$\begin{aligned}\hat{\alpha}_0 &:= \tau_{00} \pm \tau_{10} \frac{\nu_1}{\nu_0}, \\ \sigma_{j0} &:= \nu_j, & 0 \leq j \leq 2\kappa - 1, \\ \sigma_{j1} &:= \tau_{j+1,j} \sigma_{j+1,0} \pm (\tau_{jj} - \hat{\alpha}_0) \sigma_{j0} \pm \tau_{j-1,j} \sigma_{j-1,0}, & 1 \leq j \leq 2\kappa - 2.\end{aligned}$$

Then compute for  $j = 1, 2, \dots, \kappa - 1$ :

$$\begin{aligned}\hat{\beta}_j &:= \tau_{j,j-1} \frac{\sigma_{jj}}{\sigma_{j-1,j-1}}, \\ \hat{\alpha}_j &:= \tau_{jj} \pm \tau_{j+1,j} \frac{\sigma_{j+1,j} \tau_{j,j-1}}{\sigma_{jj}}, \\ \sigma_{i,j+1} &:= \tau_{i+1,i} \sigma_{i+1,j} \pm (\tau_{ii} - \hat{\alpha}_j) \sigma_{ij} \pm \tau_{i-1,i} \sigma_{i-1,j} - \hat{\beta}_j \sigma_{i,j-1}, & j \leq i < 2\kappa - j,\end{aligned}$$

where we use the property (5.3). Thus, the computation of the entries of  $H_\kappa$  requires  $2\kappa$  modified moments  $\{\nu_j\}_{j=0}^{2\kappa-1}$ . We compute these  $\nu_j$  from the residual vectors  $\{\mathbf{r}_j\}_{j=0}^{2\kappa-1}$  by (4.6). The eigenvalues of  $H_\kappa$  are estimates of eigenvalues of  $\mathbf{A}$ , and in the computed examples of §7 we computed them by the EISPACK [19] subroutine HQR. Our computational experience indicates that  $\kappa$  often should be chosen fairly small, e.g.,  $\kappa = 5$ . After each  $2\kappa - 1$  iterations by the Chebyshev method (2.1), we can determine a new matrix  $H_\kappa$  and compute its spectrum. Let  $S(H)$  denote the union of sets of eigenvalues of all the computed matrices  $H_\kappa$ . We determine the parameters  $\mathbf{d}$  and  $c$  by fitting an ellipse to the available set  $S(H)$  using the scheme outlined in §3 with  $S(\mathbf{A})$  replaced by  $S(H)$ . The spectrum of each computed matrix  $H_\kappa$  increases the set  $S(H)$ . If during the iteration with the adaptive Chebyshev method one finds that the parameters  $\mathbf{d}$  and  $c$  change insignificantly when eigenvalues of new matrices  $H_\kappa$  are included in the set  $S(H)$ , then Chebyshev iteration can typically proceed until convergence with fixed values of  $\mathbf{d}$  and  $c$  and without computing further modified moments.

## 6 Modified modified moments

This section describes an alternative way of estimating the spectrum of  $\mathbf{A}$  by using modified moments. Manteuffel [16] points out that estimating the spectrum of  $\mathbf{A}$  by the power method can give estimates biased towards eigenvalues of large magnitude. However, since eigenvalues of large magnitude are not necessarily associated with large convergence factors, they might not be the most important ones for determining good parameters  $\mathbf{d}$  and  $c$ . Manteuffel [16] proposes to apply the power method to the matrix  $\hat{\mathbf{A}}$ , defined by (2.10), to circumvent this problem; see the discussion in §2. In this section we show how eigenvalue estimates of  $\mathbf{A}$  can be computed from modified moments associated with  $\hat{\mathbf{A}}$ . These modified moments are computed by modifying the modified moments associated with  $\mathbf{A}$ . The work



and storage required for computing these modified modified moments of  $\mathbf{A}$  exceed the requirements for the scheme based on modified moments of  $\mathbf{A}$  only by a negligible amount.

The matrix  $\hat{\mathbf{A}}$  has spectral decomposition  $\hat{\mathbf{A}} = \mathbf{W}\hat{\Lambda}\mathbf{W}^{-1}$ , where

$$\hat{\Lambda} = \text{diag}[\hat{\lambda}_1, \hat{\lambda}_2, \dots, \hat{\lambda}_N] := 2g(d\mathbf{I} - \mathbf{A}). \quad (6.1)$$

Introduce the vectors

$$\hat{\mathbf{r}}_{\mathbf{k}} := p_{\mathbf{k}}(\hat{\mathbf{A}})\hat{\mathbf{r}}_0, \quad \hat{\mathbf{r}}_0 := \mathbf{r}_0, \quad (6.2)$$

where the polynomials  $p_{\mathbf{k}}$  are given by (2.3). Similarly as in §4, we introduce a complex symmetric measure  $\hat{\gamma}$  associated with  $\hat{\mathbf{A}}$  and  $\hat{\mathbf{r}}_0$ , i.e.,  $\hat{\gamma}$  has jumps of height  $\gamma_{ij} := \alpha_i \alpha_j \mathbf{w}_i^T \mathbf{w}_j$  at the points  $(\hat{\lambda}_i, \hat{\lambda}_j) \in \mathbb{C}^2$ . Then

$$\hat{\mathbf{r}}_{\mathbf{k}}^T \hat{\mathbf{r}}_l = \iint_{\mathbf{C}} p_{\mathbf{k}}(\lambda) p_l(\eta) d\hat{\gamma}(\lambda, \eta). \quad (6.3)$$

In view of (6.1), we can replace  $\hat{\gamma}$  in (6.3) by the measure  $\gamma$  introduced in §4 in the following way. Let

$$\hat{p}_{\mathbf{k}}(\lambda) := p_{\mathbf{k}}(2gd - 2g\lambda).$$

Then

$$\int_{\mathbf{C}} \int_{\mathbf{C}} p_{\mathbf{k}}(\lambda) p_l(\eta) d\hat{\gamma}(\lambda, \eta) = \int_{\mathbf{C}} \int_{\mathbf{C}} \hat{p}_{\mathbf{k}}(\lambda) \hat{p}_l(\eta) d\gamma(\lambda, \eta).$$

In particular, we can define modified moments associated with  $\hat{\mathbf{A}}$  by

$$\hat{\nu}_{\mathbf{k}} := \hat{\mathbf{r}}_{\mathbf{k}}^T \hat{\mathbf{r}}_0 = \int_{\mathbf{C}} \int_{\mathbf{C}} \hat{p}_{\mathbf{k}}(\lambda) d\gamma(\lambda, \eta). \quad (6.4)$$

We now write  $\hat{p}_{\mathbf{k}}$  as a linear combination of residual polynomials  $p_j$  of degree  $j \leq \mathbf{k}$ . Substituting this linear combination into (6.4) yields a formula for expressing the modified moment  $\hat{\nu}_{\mathbf{k}}$  in terms of modified moments  $\nu_j$  for  $j \leq \mathbf{k}$ . The  $\nu_j$  are as usual computed by (4.6). In the remainder of this section we determine the coefficients in this linear combination.

Let  $\hat{b} := -2g$  and  $\hat{a} := 2gd$ . We would like to compute constants  $c_{\mathbf{k}j}$ , such that

$$p_{\mathbf{k}}(\hat{a} + \hat{b}\lambda) = \sum_{j=0}^{\mathbf{k}} c_{\mathbf{k}j} p_j(\lambda). \quad (6.5)$$

It follows from

$$p_{\mathbf{k}}(\hat{a} + \hat{b}\lambda) = \frac{T_{\mathbf{k}}\left(\frac{d - (\hat{a} + \hat{b}\lambda)}{c}\right)}{T_{\mathbf{k}}\left(\frac{d}{c}\right)}$$

that the coefficients  $c_{\mathbf{k}j}$  in (6.5) are easily determined from the coefficients  $\hat{c}_{\mathbf{k}j}$  in

$$T_{\mathbf{k}}\left(\frac{d - (\hat{a} + \hat{b}\lambda)}{c}\right) = \sum_{j=0}^{\mathbf{k}} \hat{c}_{\mathbf{k}j} T_j\left(\frac{d - \lambda}{c}\right). \quad (6.6)$$

In order to simplify the notation, we introduce

$$\xi := \frac{d - \lambda}{c}, \quad a := \frac{d}{c}, \quad b := \hat{b} = -2g.$$

Then

$$T_k \left( \frac{d - (\hat{a} + \hat{b}\lambda)}{c} \right) = T_k(a + b\xi).$$

The coefficients  $\hat{c}_{kj}$  can be computed recursively. Assume that the  $\hat{c}_{ij}$  are known for  $0 \leq i, j \leq k$ .

Combining the three-term recurrence relation

$$T_{k+1}(a + b\xi) = 2(a + b\xi)T_k(a + b\xi) - T_{k-1}(a + b\xi)$$

with (6.6) yields

$$\sum_{j=0}^{k+1} \hat{c}_{k+1,j} T_j(\xi) = \sum_{j=0}^k 2a\hat{c}_{kj} T_j(\xi) + \sum_{j=0}^k \hat{c}_{kj} b 2\xi T_j(\xi) - \sum_{j=0}^{k-1} \hat{c}_{k-1,j} T_j(\xi) \quad (6.7)$$

$$= \sum_{j=0}^k 2a\hat{c}_{kj} T_j(\xi) + \sum_{j=1}^{k+1} b\hat{c}_{k,j-1} T_j(\xi) + \sum_{j=0}^{k-1} b\hat{c}_{k,j+1} T_j(\xi) - \sum_{j=0}^{k-1} \hat{c}_{k-1,j} T_j(\xi). \quad (6.8)$$

Identifying coefficients on the left and right hand sides of (6.7) or (6.8) gives

$$\begin{aligned} \hat{c}_{k+1,0} &= 2a\hat{c}_{k0} + b\hat{c}_{k1} - \hat{c}_{k-1,0}, \\ \hat{c}_{k+1,1} &= 2a\hat{c}_{k1} + 2b\hat{c}_{k0} + b\hat{c}_{k2} - \hat{c}_{k-1,1}, \quad k \geq 2, \\ \hat{c}_{k+1,j} &= 2a\hat{c}_{kj} + b\hat{c}_{k,j-1} + b\hat{c}_{k,j+1} - \hat{c}_{k-1,j}, \quad 2 \leq j < k, \quad k \geq 3, \\ \hat{c}_{k+1,k} &= 2a\hat{c}_{kk} + b\hat{c}_{k,k-1}, \\ \hat{c}_{k+1,k+1} &= b\hat{c}_{kk}. \end{aligned} \quad (6.9)$$

It follows from  $\mathbf{T}_0(\mathbf{X}) = 1$ ,  $T_1(\lambda) = \lambda$  and from (6.6) that

$$\hat{c}_{00} = 1, \quad \hat{c}_{10} = \frac{d}{c}, \quad \hat{c}_{11} = -2g. \quad (6.10)$$

Recalling that the residual polynomials  $p_k(\hat{A})$  satisfy

$$p_k(\hat{A}) = \frac{T_k(aI + \frac{b}{c}(dI - \mathbf{A}))}{T_k(\frac{d}{c})},$$

we obtain

$$p_k(\hat{A}) = \sum_{j=0}^k \hat{c}_{kj} \frac{T_j(\frac{d}{c})}{T_k(\frac{d}{c})} \frac{T_j(\frac{1}{c}(dI - \mathbf{A}))}{T_j(\frac{d}{c})} = \sum_{j=0}^k c_{kj} p_j(\mathbf{A}), \quad (6.11)$$

where

$$c_{kj} = \hat{c}_{kj} \frac{T_j(\frac{d}{c})}{T_k(\frac{d}{c})},$$

and the  $\hat{c}_{kj}$  can be computed recursively from (6.9) and (6.10). We remark that the coefficients  $\hat{c}_{kj}$  are real when  $k + j$  is even, and purely imaginary when  $k + j$  is odd. Therefore, since  $T_j(d/c)$  is real for  $j$  even, and purely imaginary for  $j$  odd, the coefficients  $c_{kj}$  are real for all  $k$  and  $j$ .

It follows from the definition of the modified moments and (6.11) that

$$\hat{\nu}_k = \sum_{j=0}^k c_{kj} \nu_j. \quad (6.12)$$

Formula (6.12) allows us to compute the first  $2\kappa$  modified moments associated with  $\hat{\mathbf{A}}$  from the first  $2\kappa$  modified moments associated with  $\mathbf{A}$ . Eigenvalue estimates for  $\hat{\mathbf{A}}$  can be computed from the modified moments for  $\hat{\mathbf{A}}$  and the recursion coefficients of the residual polynomials  $p_k$  in a similar manner as the eigenvalue estimates for  $\mathbf{A}$  are computed from modified moments of  $\mathbf{A}$  and recursion coefficients of the  $p_k$ . We finally note that eigenvalue estimates for  $\mathbf{A}$  are easily obtained from eigenvalue estimates for  $\hat{\mathbf{A}}$  by (6.1).

## 7 Numerical examples

This section presents numerical experiments, in which the performance of our two new adaptive Chebyshev algorithms for nonsymmetric linear systems based on modified moments are compared with the adaptive Chebyshev method based on the power methods applied to the matrices  $\mathbf{A}$  and  $\hat{\mathbf{A}}$  by Manteuffel [16] as implemented by CHEBYCODE [2], where  $\hat{\mathbf{A}}$  is defined by (2.10). All programs used are written in FORTRAN 77. Our new adaptive schemes have been implemented by using parts of CHEBYCODE [2], e.g., the subroutines for computing the Chebyshev iterates and for determining and updating the smallest ellipse containing all computed eigenvalue estimates of the matrix  $\mathbf{A}$ .

We carried out the numerical experiments on an IBM RISC 6000/550 workstation using double precision arithmetic, i.e., with approximately 15 significant digits. The test problems are derived by discretizing the elliptic partial differential equation

$$-\Delta u + 2p_1 u_x + 2p_2 u_y - p_3 u = f \quad (7.1)$$

with constant coefficients  $p_1$ ,  $p_2$  and  $p_3$  on the unit square  $\Omega := \{(x, y) : 0 \leq x, y \leq 1\}$ , and with boundary condition  $u(x, y) = 0$  on  $\partial\Omega$ . The function  $f$  is chosen so that  $u(x, y) = x e^{xy} \sin(\pi x) \sin(\pi y)$  solves (7.1). We discretize (7.1) by symmetric finite differences on a uniform  $(n + 2) \times (n + 2)$  grid, including boundary points, and use the standard five-point stencil to approximate  $\Delta u$ . This yields a linear system of  $N := n^2$  equations for  $n^2$  unknowns  $u_{ij}$ ,  $1 \leq i, j \leq n$ , where  $u_{ij}$  approximates the solution  $u$  of (7.1) at the grid point  $(ih, jh)$ ,  $h := \frac{1}{n+1}$ . We scale the linear system obtained in this

manner by  $h^2$  and write it as  $\tilde{\mathbf{A}}\mathbf{x} = \mathbf{b}$ . A typical equation of this system reads

$$(4 - p_3 h^2)u_{ij} - (1 + p_1 h)u_{i-1,j} - (1 - p_1 h)u_{i+1,j} - (1 + p_2 h)u_{i,j-1} - (1 - p_2 h)u_{i,j+1} = h^2 f_{ij},$$

where  $f_{ij} = \mathbf{f}(\mathbf{i}\mathbf{h}, \mathbf{j}\mathbf{h})$ . In order to keep the issues of interest clear, no preconditioner is used. In practical applications, however, the use of a preconditioner is often desirable. To obtain systems of equations with different properties, we modify the matrix  $\tilde{\mathbf{A}}$  by adding a multiple of the identity, i.e., we solve  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , where  $\mathbf{A} := \tilde{\mathbf{A}} + \delta\mathbf{I}$  and  $\delta \geq 0$  is a constant. As the value of  $\delta$  increases, the spectrum of the matrix  $\mathbf{A}$  is shifted away from the origin.

In the following tables “pm(A)” and “pm( $\hat{\mathbf{A}}$ )” denote the adaptive Chebyshev algorithm based on the power methods applied to the matrices  $\mathbf{A}$  and  $\hat{\mathbf{A}}$ , respectively, as implemented by the code [2]. We recall that the number of matrix-vector products required by pm(A) exceeds the number of iterations, because each fitting of the ellipse requires 4 matrix-vector products that are not used to update the iterates. The number of iterations required by the methods is listed in the column labeled “steps”. We denote our adaptive scheme based on modified moments associated with the matrix  $\mathbf{A}$  by “mm(A)”, and “mm(˜)” stands for our adaptive scheme based on modified moments associated with the matrix  $\hat{\mathbf{A}}$ . The column in the tables labeled “maxadapt” shows the maximum number of times the ellipse is fitted in the schemes mm(A) and mm( $\hat{\mathbf{A}}$ ). The column labeled “frequency” show how often the ellipse is fitted. For instance, if maxadapt = 10 and frequency = 20, then the ellipse is fitted after every 20 iterations until the ellipse has been fitted 10 times. In all the examples, we choose  $\kappa = 5$ , i.e., the computed eigenvalue estimates of  $\mathbf{A}$  are eigenvalues of  $5 \times 5$  tridiagonal matrices. After each fitting of the ellipse only  $2\kappa = 10$  modified moments (4.6) have to be computed, and therefore only 10 inner products are computed, independently of the frequency  $\geq 10$  chosen. Moreover, after the ellipse has been fitted maxadapt times no more modified moments, and therefore no more inner products, are computed. We remark that our code for the adaptive Chebyshev algorithm based on modified moments is a research code and lacks the sophistication of a production code. We believe that its performance can be improved by careful coding and by implementation of strategies for choosing the frequency and maxadapt parameters.

In the derivation of the modified Chebyshev algorithm of Section 5 we assumed that the tridiagonal matrix  $\mathbf{H}$  exists. If during the computations it would turn out that the  $\kappa \times \kappa$  leading principal submatrix  $\mathbf{H}_\kappa$  of  $\mathbf{H}$  does not exist, then the modified Chebyshev algorithm is curtailed and a submatrix of  $\mathbf{H}_\kappa$  is determined, whose spectrum yields eigenvalue estimates of  $\mathbf{A}$ .

Example 7.1. We select  $p_1 = 60$ ,  $p_2 = 80$ ,  $p_3 = 40$  and  $\delta = 0.05$ . Table 7.1 shows that the adaptive

adaptive method	maxadapt	frequency	steps	$\ r_{last}\ /\ r_0\ $
pm(A)			276	.59D-10
mm(A)	15	10	275	.60D-10
mm(A)	10	20	242	.48D-10
mm(A)	6	30	235	.56D-10
mm(A)	7	35	229	.60D-10
mm(A)	7	40	229	.60D-10

Table 7.1:  $p_1 = 60$ ,  $p_2 = 80$ ,  $p_3 = 40$ ;  $\delta = 0.05$ ;  $N = 10,000$

adaptive method	maxadapt	frequency	steps	$\ r_{last}\ /\ r_0\ $
pm(A)			271	.63D-10
mm(A)	4	40	224	.32D-10
mm(A)	4	45	224	.32D-10

Table 7.2:  $p_1 = 60$ ,  $p_2 = 80$ ,  $p_3 = 40$ ;  $\delta = 0.05$ ;  $N = 25,600$

schemes mm(A) and pm(A) achieve convergence in approximately the same number of Chebyshev iterations when the order of  $\mathbf{A}$  is  $N = 10,000$  and the adaptive procedure in mm(A) is called every 10 iterations. It is clear from Table 7.1 that, as we reduce the frequency of calls to the adaptive procedure, the number of iterations necessary to achieve roughly the same residual error is reduced by up to 17%. This depends on that Chebyshev iteration is restarted after each fitting of the ellipse. Tables 7.2-7.3 show that a similar decrease in the number of iterations is obtained for larger systems as well. This example, as well as many of the following ones, illustrates that a careful implementation of the scheme pm(A) should include strategies for choosing the parameters frequency and maxadapt. cl

Example 7.2. We select  $p_1 = 60$ ,  $p_2 = 80$ ,  $p_3 = 40$  and  $\delta = 0.02$ . We remark that with this choice of  $\delta$  the spectrum of  $\mathbf{A}$  is closer to the origin than in Example 7.1. Table 7.4 shows that, as the spectrum of the matrix  $\mathbf{A}$  of order  $N = 10,000$  is moved closer to the origin, our adaptive scheme mm(A) requires significantly fewer iterations than the scheme pm(A). Similar behavior can be seen when  $N = 40,000$ ;

adaptive method	maxadapt	frequency	steps	$\ r_{last}\ /\ r_0\ $
pm(A)			270	.60D-10
mm(A)	9	30	278	.60D-10
mm(A)	9	25	229	.24D-10
mm(A)	9	20	231	.42D-10

Table 7.3:  $p_1 = 60$ ,  $p_2 = 80$ ,  $p_3 = 40$ ;  $\delta = 0.05$ ;  $N = 40,000$

adaptive method	maxadapt	frequency	steps	$\ r_{last}\ /\ r_0\ $
pm(A)			412	.42D-10
pm(A)			292	.38D-10
mm(A)	9	20	301	.37D-10
mm(A)	9	30	286	.32D-10
mm(A)	9	40	286	.32D-10

Table 7.4:  $p_1 = 60, p_2 = 80, p_3 = 40; \delta = 0.02; N = 10,000$

adaptive method	maxadapt	frequency	steps	$\ r_{last}\ /\ r_0\ $
pm(A)			408	.38D-10
pm(A)			352	.35D-10
mm(A)	9	25	321	.12D-10
mm(A)	9	30	306	.14D-10

Table 7.5:  $p_1 = 60, p_2 = 80, p_3 = 40; \delta = 0.02; N = 40,000$

see Table 7.5. We have observed that when the matrix  $A$  has eigenvalues very close to the origin, the scheme pm(A) often requires fewer iterations than pm(A). For large such systems, our schemes mm(A) and mm(A) typically require even fewer iterations. Table 7.6 provides another illustration of this performance. In this table method pm(A) requires 67% more iterations and method pm(A) 39% more iterations than the scheme mm(A).  $\square$

Example 7.3. We select  $p_1 = 60, p_2 = 80, p_3 = 40$  and  $\delta = 0.01$ . When  $N = 10,000$  the scheme pm(A) yields a much larger error after 1000 iterations than mm(A) after only 647 iterations; see Table 7.7. The dominating work required for determining eigenvalue estimates by method mm(A) is the computation of 90 inner products (4.6). As the size of the linear system increases the scheme mm(A) performs significantly better than pm(A) and pm(A). Table 7.8 shows that when  $N = 40,000$ , the latter schemes require at least 42% more iterations than mm(A) to achieve a comparable reduction in the norm of the residual vector.  $\square$

adaptive method	maxadapt	frequency	steps	$\ r_{last}\ /\ r_0\ $
pm(A)			902	.19D-12
pm(A)			751	.18D-12
mm(A)	9	35	540	.19D-12

Table 7.6:  $p_1 = 80, p_2 = 80, p_3 = 40; \delta = 0.015; N = 40,000$

adaptive method	maxadapt	frequency	steps	$\ \mathbf{r}_{last}\ /\ \mathbf{r}_0\ $
pm(A)			1000	.16D-10
mm(A)	9	25	674	.13D-12
mm(A)	9	30	647	.13D-12
mm(A)	9	35	647	.13D-12

Table 7.7:  $p_1 = 60, p_2 = 80, p_3 = 40; \delta = 0.01; N = 10,000$

adaptive method	maxadapt	frequency	steps	$\ \mathbf{r}_{last}\ /\ \mathbf{r}_0\ $
pm(A)			1000	.17D-10
pm(A)			983	.12D-12
mm(A)	9	35	694	.13D-12

Table 7.8:  $p_1 = 60, p_2 = 80, p_3 = 40; \delta = 0.01; N = 40,000$

Example 7.4. We select  $p_1 = 30, p_2 = 40, p_3 = 40$  and  $\delta = 0$ . Let  $N = 2500$ . This example illustrates that the scheme mm(A) can give faster convergence than the method mm(A). Moreover, both methods mm(A) and mm(A) converge, while the methods pm(A) and pm(A) do not; see Table 7.9. We remark that in most of the previous examples the methods mm(A) and mm(A) display about the same rate of convergence. cl

## 8 Conclusions

This paper presents two adaptive Chebyshev algorithms for solving large, sparse nonsymmetric linear systems based on modified moments. A major advantage of these scheme is that they require fewer N-vectors be stored in computer memory than adaptive schemes based on the power method. Moreover, our numerical examples illustrate that the schemes based on modified moments often yield significantly faster convergence if the matrix has eigenvalues close to the origin. When the eigenvalues are not close to the origin, adaptive schemes based on the power method can yield as rapid convergence

adaptive method	maxadapt	frequency	steps	$\ \mathbf{r}_{last}\ /\ \mathbf{r}_0\ $
pm(A)			1000	.43
pm(A)			1000	.39
mm(A)	$\infty$	10	1000	.82D-2
mm(A)	$\infty$	10	1000	.52D-4

Table 7.9:  $p_1 = 30, p_2 = 40, p_3 = 40; \delta = 0; N = 2,500$

as our schemes based on modified moments. The choices of how often the ellipse is to be fitted is important for the performance of our new methods. Computed examples indicate that the ellipse does not have to be fitted many times, and, therefore, the iterations can be carried out by evaluating only fairly few inner products. Our scheme therefore is attractive for implementation on parallel MIMD and SIMD computers. Such implementations should adaptively determine how frequently the ellipse ought to be fitted, and when an ellipse has been determined that can be used until convergence.

**Acknowledgement** We would like to thank Steve Ashby for providing the code described in [2].

## References

- [1] W.E. Arnoldi, The principle of minimized iterations in the solution of the matrix eigenvalue problem, *Quart. Appl. Math.*, **9** (1951), pp. 17-29.
- [2] S.F. Ashby, CHEBYCODE: a FORTRAN implementation of Manteuffel's adaptive Chebyshev algorithm, Report UIUCDCS-R-85-1203, Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, IL, 1985.
- [3] R.G. Bartle, *Elements of Real Analysis*, Wiley, New York, 1964.
- [4] F. Chatelin and S. Godet-Thobie, Stability analysis in aeronautical industries, in *High Performance Computing II*, eds. M. Durand and F. El Dabaghi, Elsevier Science Publishers, 1991, pp. 415-422.
- [5] J.J. Dongarra, I.S. Duff, D.C. Sorensen and H.A. van der Vorst, Solving *Linear Systems on Vector and Shared Memory Computers*, SIAM, Philadelphia, 1991.
- [6] M. Eiermann, On semiiterative methods generated by Faber polynomials, *Numer. Math.*, **56** (1989), pp. 139-156.
- [7] H.C. Elman, Y. Saad and P.E. Saylor, A hybrid Chebyshev Krylov subspace algorithm for solving nonsymmetric systems of linear equations, *SIAM J. Sci. Stat. Comput.*, **7** (1986), pp. 840-855.
- [8] B. Fischer and R. Freund, Chebyshev polynomials are not always optimal, *J. Approx. Theory*, **65** (1991), pp. 261-272.
- [9] W. Gautschi, Construction of Gauss-Christoffel quadrature formulas, *Math. Comp.*, **22** (1968), pp. 251-270.
- [10] W. Gautschi, On generating orthogonal polynomials, *SIAM J. Sci. Stat. Comput.*, **3** (1982), pp. 289-317.
- [11] G.H. Golub and M.H. Gutknecht, Modified moments for indefinite weight functions, *Numer. Math.*, **57** (1990), pp. 607-624.
- [12] G.H. Golub and M. Kent, Estimates of eigenvalues for iterative methods, *Math. Comp.*, **53** (1989), pp. 619-626.
- [13] G.H. Golub and R.S. Varga, Chebyshev semi-iterative methods, successive over-relaxation methods, and second order Richardson iterative methods I+II, *Numer. Math.*, **3** (1961), pp. 147-168.



- [14] D. Ho, Tchebyshev acceleration technique for large scale nonsymmetric matrices, *Numer. Math.*, **56** (1990), pp. 721-734.
- [15] T.A. Manteuffel, The Chebyshev iteration for nonsymmetric linear systems, *Numer. Math.*, **28** (1977), pp. 307-327.
- [16] T.A. Manteuffel, Adaptive procedure for estimation of parameters for the nonsymmetric Chebyshev iteration, *Numer. Math.*, **31** (1978), pp. 187-208.
- [17] N.M. Nachtigal, L. Reichel and L.N. Trefethen, A hybrid GMRES algorithm for nonsymmetric linear systems, *SIAM J. Matrix Anal. Appl.*, **13** (1992), to appear.
- [18] R.A. Sack and A.F. Donovan, An algorithm for Gaussian quadrature given modified moments, *Numer. Math.*, **18** (1972), pp. 465-478.
- [19] B.T. Smith, J.M. Boyle, Y. Ikebe, V.C. Klema and C.B. Moler, *Matrix Eigensystem Routines: EISPACK Guide*, 2nd ed., Springer-Verlag, New York, NY, 1970.
- [20] R.S. Varga, *Matrix Iterative Analysis*, Prentice-Hall, Englewood Cliffs, NJ, 1962.
- [21] J.C. Wheeler, Modified moments and Gaussian quadratures, *Rocky Mountain J. Math.*, **4** (1974), pp. 287-295.

