

# NEW BAND TOEPLITZ PRECONDITIONERS FOR ILL-CONDITIONED SYMMETRIC POSITIVE DEFINITE TOEPLITZ SYSTEMS

D. NOUTSOS\* AND P. VASSALOS†

**Abstract.** It is well known that Preconditioned Conjugate Gradient (PCG) methods are widely used to solve ill-conditioned Toeplitz linear systems  $T_n(f)x = b$ . In this paper we present a new preconditioning technique for the solution of symmetric Toeplitz systems generated by nonnegative functions  $f$  with zeros of even order. More specifically,  $f$  is divided by the appropriate trigonometric polynomial  $g$  of the smallest degree, with zeros the zeros of  $f$ , to eliminate its zeros. Using rational approximation we approximate  $\sqrt{\frac{f}{g}}$  by  $\frac{p}{q}$  and consider  $\frac{p^2g}{q^2}$  as a very satisfactory approximation of  $f$ . We propose the matrix  $M_n = B_n^{-1}(p)B_n(p^2g)B_n^{-1}(p)$  as a preconditioner whence a good clustering of the spectrum of its preconditioned matrix is obtained. We also show that the proposed technique can be very flexible, a fact that is confirmed by various numerical experiments so that in many cases it constitutes a much more efficient strategy than the existing ones.

**Key words.** low rank correction, Toeplitz matrix, conjugate gradient, rational interpolation and approximation, preconditioner

**AMS subject classifications.** Primary 65F10, 65F15

**1. Introduction.** In this paper we use and analyze band Toeplitz matrices as preconditioners for the solution of the  $n \times n$  ill-conditioned symmetric and positive definite Toeplitz system

$$(1.1) \quad T_n(f)x = b$$

by the Preconditioned Conjugate Gradient (PCG) method, where the matrix  $T_n(f) \in \mathbb{R}^{n \times n}$  is produced by a real-valued, even,  $2\pi$ -periodic function defined in the fundamental interval  $[-\pi, \pi]$ . Then, the  $(j, k)$  element of  $T_n(f)$  is given by the Fourier coefficient of  $f$ , i.e

$$T_n(f)_{j,k} = T_{j-k} = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) e^{-i(j-k)x} dx, \quad 1 \leq j, k < n,$$

where  $i$  is the imaginary unit.

Toeplitz matrices arise very often in a wide variety of applications, as e.g., in the numerical solution of differential equations using finite differences, in statistical problems (linear prediction), in Wiener-Hopf kernels, in Markov chains, in image and signal processing, e.t.c. (see [8], [3], [19]). The generating function  $f$  plays a significant role in the location and distribution of the eigenvalues of Toeplitz matrix [8], [4] and in many cases is a priori known. As it is known for the spectrum of  $T_n(f)$  there holds  $\sigma(T_n(f)) \subseteq [\text{ess inf } f, \text{ess sup } f]$ .

Superfast direct methods can solve system (1.1) in  $O(n \log^2 n)$  operations, but their stability properties for ill-conditioned Toeplitz matrices are still unclear; see, for instance, [3].

---

\*Department of Mathematics, University of Ioannina, GR-451 10, Ioannina, Greece (dnoutsos@cc.uoi.gr).

†Department of Mathematics, University of Ioannina, GR-451 10, Ioannina, Greece (pvassal@pythagoras.math.uoi.gr). The research of this author was supported by Hellenic Foundation of Scholarships (HFS).

The classical iterative methods such as Jacobi, Gauss-Seidel and SOR are not effective since the associated spectral radius tends to 1 for large  $n$ . The method which is widely used for the solution of such systems is the PCG method. The factors that affect the convergence features of this method are the magnitude of the condition number  $\kappa_2(T_n(f))$  and the distribution of the eigenvalues. So a good preconditioner must cluster the eigenvalues of the preconditioned system as much as possible and make the eigenvalues that might lie outside the cluster have magnitude independent of  $n$ .

If the generating function is continuous and positive then problem (1.1) will not be ill-conditioned and the condition number can not increase proportionally to  $n$  although it can be very large. In this case system (1.1) can be handled by using a preconditioner belonging to some Trigonometric matrix algebras (circulant,  $\tau$ , Hartley, [18], [17], [9]) or by band Toeplitz preconditioners with weakly increasing bandwidth defined by a polynomial operator  $\mathcal{S}_n$  as was proposed in [16]. Theoretically, the latter class of preconditioners seems to perform better as  $n \rightarrow \infty$  since the number of PCG iterations tends to 1 while in the former cases this number tends to a constant.

When  $f$  has any zeros, then system (1.1) is ill-conditioned and the condition number  $\kappa_2(T_n(f))$  increases proportionally to  $n^\alpha$  where  $\alpha$  is the largest number of the multiplicities of the zeros of  $f$  [4], [14]. To best handle this case it is necessary to know the number of the zeros of  $f$ . If this number is not even then the most suitable technique for this situation [13], fails to make the condition number of the preconditioned matrix independent of its dimension  $n$  and the problem is still open. On the other hand things dramatically change when the number of zeros is even.

In this case, it was R. Chan [4] who first proposed as a preconditioner for system (1.1) the Toeplitz band matrix  $B_n(g)$  whose generating function  $g$  is a trigonometric polynomial that has the same zeros with the same multiplicities as those of  $f$ . Next, in [5],  $g$  was not only considered as having the zeros of  $f$  but also its degree was increased so that it provided additional degrees of freedom to approximate  $f$  and to minimize the relative error  $\|\frac{f-g}{g}\|_\infty$  over all trigonometric polynomials  $g$  of a fixed degree  $l$ . The generating function  $g$  is then computed by the Remez algorithm, which can be very expensive, from the computational point of view, especially when  $f$  has a large number of zeros.

Recently, Serra [15] has extended this method by proposing alternative techniques to minimize  $\|\frac{f-g}{g}\|_\infty$ . More specifically, he chose as  $g$ ,  $z_k g_{l-k}$  where  $z_k$  is the trigonometric polynomial of minimum degree  $k$  that has all the zeros of  $f$  with their multiplicities and  $g_{l-k}$  is the trigonometric polynomial of degree  $l-k$  which is the best Chebyshev approximation of  $\hat{f} = \frac{f}{z_k}$  from the space  $\mathcal{P}_{l-k}$  of all trigonometric polynomials of degree at most  $l-k$ . In addition, in the same work [15], it was also proposed another way of constructing  $g_{l-k}$  by interpolating  $\hat{f}$  at the  $l-k+1$  zeros of the  $(l-k+1)$ -st degree Chebyshev polynomial of the first kind.

We remark that it has been proved [7], that preconditioners belonging to the aforementioned matrix algebra, when they can be defined, produce weak clustering, i.e., the eigenvalues of the preconditioned matrix are such that for every  $\epsilon > 0$  there exists a positive  $\beta$  so that, except for rare exceptions,  $O(n^\beta)$  of the eigenvalues lie in the interval  $(0, \epsilon)$ .

In this paper we extend the previous methods in order to achieve a better clustering for the eigenvalues of the preconditioned matrix and propose a way of constructing a class of preconditioners based on rational approximation or on interpolation to the positive and continuous function  $\sqrt{\frac{f}{z_k}}$  with  $z_k$  defined previously.

The outline of the present work is as follows. In Section 2 we recall some useful issues about the rational approximation, while in Section 3 we introduce the technique of constructing the new class of preconditioners based on rational approximation to  $\sqrt{\frac{f}{z_\rho}}$  and analyze the convergence of the PCG method. In Section 4 we study the flexibility and possible modifications of our method, analyze its cost per iteration and compare it with that of previous techniques. Finally, in Section 5, results of illustrative numerical experiments are exhibited and concluding remarks are made.

**2. Preliminaries.** In what follows we assume that the generating function  $f$  is defined in  $[-\pi, \pi]$ , is  $2\pi$ -periodic, continuous, nonnegative and has zeros of even order.

We define by  $z_k$  the trigonometric polynomial of minimum degree  $k$  containing all the zeros of  $f$  with their multiplicities. Then we define  $r_{lm} = \frac{p_l}{q_m}$  as the best rational approximation of  $\hat{f} = \sqrt{\frac{f}{z_k}}$  in the uniform norm, i.e.,

$$\|\hat{f} - r_{lm}\|_\infty = \min_{r \in \mathcal{R}(l, m)} \|\hat{f} - r\|_\infty,$$

where  $\mathcal{R}(l, m)$  denotes the set of rational functions  $r$ , with  $p \in \mathcal{P}_l$ ,  $q \in \mathcal{P}_m$  and  $r$  is irreducible, that is  $p$  and  $q$  have no zeros in common.

It is known that when  $f$  belongs to some special class of functions [10] then the order of magnitude of the maximum error of an approximation from the space  $\mathcal{R}(l, m)$  is better than the corresponding error in the space  $\mathcal{P}(l + m)$ . In general, we hope that taking advantage of the flexible nature of rational functions this set will be a stronger tool than its competitor the polynomial one. For example, it is obvious that polynomials are not suitable for approximating functions having sharp peaks near the center of their ranges and are slowly varying when  $|x|$  increases. Such kind of behavior can be obtained by continuous functions which are not differentiable at some points. However, it is easy to overcome this difficulty by using rational functions.

The next theorem establishes the fact that rational approximation of continuous functions in  $[-\pi, \pi]$  is always possible and unique.

**THEOREM 2.1.** *Let  $f$  in  $C[-\pi, \pi]$ . Then there exists  $r^* \in \mathcal{R}(l, m)$  such that*

$$\|f - r^*\| < \|f - r\|$$

for all  $r \in \mathcal{R}(l, m)$ ,  $r \neq r^*$ .

*Proof.* See [12], pp. 121, 125.  $\square$

**3. Construction of the Preconditioner.** Let  $f$  be a  $2\pi$ -periodic, nonnegative function belonging to  $C[-\pi, \pi]$  with zeros  $x_1, x_2, \dots, x_s$  of multiplicities  $2\mu_1, 2\mu_2, \dots, 2\mu_s$ , respectively, and  $2\mu_1 + 2\mu_2 + \dots + 2\mu_s = \rho$ . First, we define

$$z_\rho = \prod_{i=1}^s (1 - \cos(x - x_i))^{\mu_i}$$

which is the trigonometric polynomial of minimum degree  $\rho$  having all the zeros of  $f$ . By dividing  $f$  by  $z_\rho$ , all its zeros are eliminated and the ratio  $\frac{f}{z_\rho}$  becomes a real positive function.

Then, we define the function  $\hat{f} = \sqrt{\frac{f}{z_\rho}}$  and approximate it with the rational trigonometric function  $r_{l, m} = \frac{p_l}{q_m}$  where  $l, m$  are the degrees of the numerator and the

denominator, respectively. Since  $\frac{p_l}{q_m}$  is the best rational approximation of  $\sqrt{\frac{f}{z_\rho}}$  for certain  $l$  and  $m$  we are led to the conclusion that  $\frac{p_l^2}{q_m^2}$  may be a good approximation of  $\frac{f}{z_\rho}$ . This means that there exists a small  $\epsilon > 0$  such that

$$\left\| \frac{f}{z_\rho} - \frac{p_l^2}{q_m^2} \right\|_\infty < \epsilon$$

or, equivalently, that there exists a small  $\delta > 0$  such that

$$\left\| \frac{q_m^2}{z_\rho p_l^2} f - 1 \right\|_\infty < \delta.$$

The last inequality means that the values of  $\frac{q_m^2}{z_\rho p_l^2} f$  are clustered in a small region near the constant number 1. In matrix analog, this means that taking  $T_n \left( \frac{z_\rho p_l^2}{q_m^2} \right)$  as a preconditioner matrix for the solution of (1.1), the eigenvalues of  $T_n^{-1} \left( \frac{z_\rho p_l^2}{q_m^2} \right) T_n(f)$  are clustered in a small region near 1 and the PCG method will become very fast. Unfortunately, this matrix is a full Toeplitz matrix, is hard to construct, is costly to invert and so it is useless as a preconditioner. Instead, we are led to the idea of separating the numerator and the denominator of the ratio  $\frac{z_\rho p_l^2}{q_m^2}$  and use as a preconditioner matrix the product of three band Toeplitz matrices. More specifically, the preconditioner we propose for the solution of system (1.1) is

$$(3.1) \quad M_n = B_{nm}^{-1}(q) B_{n\hat{l}}(p^2 z_\rho) B_{nm}^{-1}(q), \quad \hat{l} = 2l + \rho,$$

where the second index in the matrices represents their halfbandwidth, while the first one their dimension. The following statements prove the basic assumptions a preconditioner must satisfy and also describe the spectrum of the preconditioned matrix  $M_n^{-1} T_n$ .

**THEOREM 3.1.** *The matrix  $M_n$  is symmetric and positive definite for every  $n$ .*

*Proof.* Its symmetry is implied directly from the definition (3.1). On the other hand, the eigenvalues of  $B_{n\hat{l}}(p^2 z_\rho)$  belong to the interval  $(\min p_l^2 z_\rho, \max p_l^2 z_\rho)$ , where  $0 = \min p_l^2 z_\rho < \max p_l^2 z_\rho \leq 2^\rho$ . Therefore,  $B_{n\hat{l}}(p^2 z_\rho)$  is symmetric and positive definite. Furthermore,  $q_m$  has no zeros in  $[-\pi, \pi]$  because it results from the rational approximation to a function which is strictly positive in  $[-\pi, \pi]$ . So,  $B_{nm}(q)$  is symmetric and invertible. Then, for every  $x \in \mathbb{R}^n$ ,  $x \neq 0$ , we have

$$x^T M_n x = x^T B_{nm}^{-1}(q) B_{n\hat{l}}(p^2 z_\rho) B_{nm}^{-1}(q) x = y^T B_{n\hat{l}}(p^2 z_\rho) y > 0,$$

where  $y = B_{nm}^{-1}(q)x$ . Hence  $M_n$  is symmetric and positive definite.  $\square$

Theorem 3.1 suggests that the matrix  $M_n$  can be taken as a preconditioner matrix. It then remains to study the convergence rate of the PCG method or, equivalently, how the eigenvalues of the matrix  $M_n^{-1} T_n$  are distributed. For this, we give without proof the following Lemma and then we state and prove our main result in Theorem 3.2.

**LEMMA 3.1.** *Suppose  $A, B \in \mathbb{R}^{n \times n}$  are symmetric matrices such that*

$$A = B + \epsilon c c^T,$$

where  $c \in \mathbb{R}^n$ ,  $c^T c = 1$ . If  $\epsilon > 0$  then

$$\lambda_1(B) \leq \lambda_1(A) \leq \lambda_2(B) \leq \cdots \leq \lambda_n(B) \leq \lambda_n(A)$$

while if  $\epsilon \leq 0$ , then

$$\lambda_1(A) \leq \lambda_1(B) \leq \lambda_2(A) \leq \cdots \leq \lambda_n(A) \leq \lambda_n(B)$$

provided that the eigenvalues are labeled in nondecreasing order of magnitude. In either case

$$\lambda_k(A) = \lambda_k(B) + t_k \epsilon, \quad k = 1, 2, \dots, n,$$

where  $t_k \geq 0$ ,  $k = 1, 2, \dots, n$ , and  $\sum_{k=1}^n t_k = 1$ .

*Proof.* See Wilkinson [20], pp. 97-98.  $\square$

**THEOREM 3.2.** Let  $\lambda_i(M_n^{-1}T_n)$ ,  $i = 1(1)n$ , denote the eigenvalues of  $M_n^{-1}T_n$  and  $m$  the degree of the denominator  $q_m$  of the rational approximation. Then, at least  $n - 4m$  eigenvalues of the preconditioned matrix lie in  $(h_{\min}, h_{\max})$ , at most  $2m$  are greater than  $h_{\max}$  and at most  $2m$  are in  $(0, h_{\min})$ , where  $h = \frac{f q^2}{p^2 z_\rho}$ .

*Proof.* Obviously the matrix

$$M_n^{-1}T_n = B_{nm}(q)B_{ni}^{-1}(p^2 z_\rho)B_{nm}(q)T_n(f)$$

is similar to the matrix

$$(3.2) \quad B_{ni}^{-\frac{1}{2}}(p^2 z_\rho)B_{nm}(q)T_n(f)B_{nm}(q)B_{ni}^{-\frac{1}{2}}(p^2 z_\rho).$$

Then, since  $B_{nm}(q)$  is a band matrix with halfbandwidth  $m$ , the matrix

$$B_{nm}(q)T_n(f)B_{nm}(q)$$

can be written as a sum of a Toeplitz matrix and a low rank correction matrix, i.e.,

$$(3.3) \quad B_{nm}(q)T_n(f)B_{nm}(q) = T_n(q^2 f) + \Delta,$$

where  $\Delta$  is a symmetric 'border' matrix with nonzero elements only in the first and last  $m$  rows and columns. So  $\text{rank}(\Delta) \leq 4m$  is independent of  $n$ . Then, from (3.2) and (3.3) we obtain that

$$(3.4) \quad \overbrace{B_{ni}^{-\frac{1}{2}}(p^2 z_\rho)B_{nm}(q)T_n(f)B_{nm}(q)B_{ni}^{-\frac{1}{2}}(p^2 z_\rho)}^E = \overbrace{B_{ni}^{-\frac{1}{2}}(p^2 z_\rho)T_n(q^2 f)B_{ni}^{-\frac{1}{2}}(p^2 z_\rho)}^{\tilde{E}} + B_{ni}^{-\frac{1}{2}}(p^2 z_\rho)\Delta B_{ni}^{-\frac{1}{2}}(p^2 z_\rho).$$

Since a matrix product does not have rank larger than that of each of the factors involved, there exist  $\alpha_i > 0$ ,  $c_i \in \mathbb{R}^n$ ,  $i = 1(1)m_+$ , and  $\beta_i > 0$ ,  $d_i \in \mathbb{R}^n$ ,  $i = 1(1)m_-$ , with  $m_+ + m_- \leq 4m$ , such that (3.4) can be written as

$$E - \tilde{E} = \sum_{i=1}^{m_+} \alpha_i c_i c_i^T - \sum_{i=1}^{m_-} \beta_i d_i d_i^T.$$

So applying successively  $m_+ + m_-$  times Lemma 3.1 gives

$$h_{\min} \leq \lambda_i(E) \leq h_{\max}, \quad m_- < i \leq n - m_+,$$

and the theorem is proved.  $\square$

It is clear from the previous analysis and statements that contrary to what happens with other band Toeplitz preconditioners, the one we propose of the 'premultiplier' matrix  $B_{nm}(q)$ , may make some of the eigenvalues lie outside the approximation interval  $[h_{\min}, h_{\max}]$ . We will prove now that the spectral radius of the preconditioned matrix is bounded by a constant number independent of  $n$ . For this, first, we state and prove the following lemma.

**LEMMA 3.2.** *Let  $B_n$  be a  $n \times n$  symmetric and positive definite band Toeplitz matrix with halfbandwidth  $s$ . Then the  $k \times k$  principal and trailing submatrices of  $B_n^{-1}$  as well as the  $k \times k$  submatrices consisting from the first  $k$  rows and the last  $k$  columns (right upper corner) or from the last  $k$  rows and the first  $k$  columns (left lower corner) of  $B_n^{-1}$ , are bounded for every fixed  $k$  independent of  $n$ .*

*Proof.* For principal and trailing submatrices, this property has been proved in [6] for  $k = s$ . We will prove the validity of this property for  $k = s + 1$  and the proof of every fixed  $k$  can be completed by induction. From the fundamental relation

$$\sum_{l=1}^{s+1} b_{1l}(B_n^{-1})_{lj} = \delta_{1j},$$

where  $\delta_{1j}$  is the Kroneker  $\delta$ , we obtain successively that

$$(3.5) \quad (B_n^{-1})_{s+1,j} = \frac{1}{b_{1,s+1}} \left( \delta_{1j} - \sum_{l=1}^s b_{1l}(B_n^{-1})_{lj} \right), \quad j = 1, 2, \dots, s.$$

Since all the elements in the righthand side of (3.5) are bounded, so are the elements  $(B_n^{-1})_{s+1,j}$ ,  $j = 1, 2, \dots, s$ . From the symmetry of  $B_n^{-1}$  we obtain that the elements  $(B_n^{-1})_{j,s+1}$ ,  $j = 1, 2, \dots, s$ , are also bounded. One more application of (3.5) for  $j = s + 1$ , gives us that the element  $(B_n^{-1})_{s+1,s+1}$  is bounded and the proof for the principal submatrices is complete. Since,  $B_n^{-1}$  is a persymmetric matrix the elements of the trailing matrix are the same as those of the principal one in reverse order. So the  $k \times k$  trailing matrix is also bounded.

It remains to prove the validity of the property for the submatrices in the right upper corner and in the left lower corner of  $B_n^{-1}$ . These matrices are transposes of each other due to the symmetry of  $B_n^{-1}$ . From the positive definiteness of  $B_n^{-1}$  we have that

$$|(B_n^{-1})_{ij}| < \frac{(B_n^{-1})_{ii} + (B_n^{-1})_{jj}}{2}, \quad i = 1, \dots, k, \quad j = n - k + 1, \dots, n.$$

The elements in the righthand side are the diagonal elements of the  $k \times k$  principal and trailing submatrices, respectively, which are bounded and the proof is complete.  $\square$

The following theorem proves that the eigenvalues of  $M^{-1}T$  have an upper bound.

**THEOREM 3.3.** *Under the assumptions of Theorem 3.2 there exists a constant  $c$ , independent of  $n$ , such that  $\rho(M_n^{-1}T_n(f)) \leq c$ , for every  $n$ .*

*Proof.* We begin the proof by using some relations connecting the spectral radii and the Rayleigh quotients of symmetric matrices. The fact that all the matrices are

positive definite, is also used.

$$\begin{aligned}
 \rho(M_n^{-1}T_n(f)) &= \rho\left(B_{nm}(q)B_{ni}^{-1}(p^2z_\rho)B_{nm}(q)T_n(f)\right) \\
 &= \rho\left(B_{ni}^{-\frac{1}{2}}(p^2z_\rho)B_{nm}(q)T_n(f)B_{nm}(q)B_{ni}^{-\frac{1}{2}}(p^2z_\rho)\right) \\
 &= \max_{x \neq 0} \frac{x^T B_{ni}^{-\frac{1}{2}}(p^2z_\rho)B_{nm}(q)T_n(f)B_{nm}(q)B_{ni}^{-\frac{1}{2}}(p^2z_\rho)x}{x^T x} \\
 &= \max_{x \neq 0} \left( \frac{x^T T_n(f)x}{x^T B_{nm}^{-1}(q)B_{ni}(p^2z_\rho)B_{nm}^{-1}(q)x} \cdot \frac{x^T B_{ni}(p^2z_\rho)x}{x^T B_{ni}(p^2z_\rho)x} \right) \\
 (3.6) \quad &= \max_{x \neq 0} \left( \frac{x^T T_n(f)x}{x^T B_{ni}(p^2z_\rho)x} \cdot \frac{x^T B_{ni}(p^2z_\rho)x}{x^T B_{nm}^{-1}(q)B_{ni}(p^2z_\rho)B_{nm}^{-1}(q)x} \right) \\
 &\leq \max_{x \neq 0} \frac{x^T T_n(f)x}{x^T B_{ni}(p^2z_\rho)x} \cdot \max_{x \neq 0} \frac{x^T B_{ni}(p^2z_\rho)x}{x^T B_{nm}^{-1}(q)B_{ni}(p^2z_\rho)B_{nm}^{-1}(q)x} \\
 &= M_1 \max_{x \neq 0} \frac{x^T B_{nm}(q)B_{ni}(p^2z_\rho)B_{nm}(q)x}{x^T B_{ni}(p^2z_\rho)x} \\
 &= M_1 \max_{x \neq 0} \frac{x^T \left( B_{n,i+2m}(q^2p^2z_\rho) + \Delta \right) x}{x^T B_{ni}(p^2z_\rho)x} \\
 &\leq M_1 \left( M_2 + \max_{x \neq 0} \frac{x^T \Delta x}{x^T B_{ni}(p^2z_\rho)x} \right) \\
 &= M_1 \left( M_2 + \rho \left( B_{ni}^{-1}(p^2z_\rho) \Delta \right) \right).
 \end{aligned}$$

In (3.6) we have taken  $M_1 = \max_{x \neq 0} \frac{x^T T_n(f)x}{x^T B_{ni}(p^2z_\rho)x} = \rho \left( B_{ni}^{-1}(p^2z_\rho)T_n(f) \right)$  and  $M_2 = \max_{x \neq 0} \frac{x^T B_{n,i+2m}(q^2p^2z_\rho)x}{x^T B_{ni}(p^2z_\rho)x} = \rho \left( B_{ni}^{-1}(p^2z_\rho)B_{n,i+2m}(q^2p^2z_\rho) \right)$  which are bounded, since the generating functions  $\frac{f}{p^2z_\rho}$  and  $\frac{q^2p^2z_\rho}{p^2z_\rho} = q^2$ , respectively, are bounded functions in  $[-\pi, \pi]$ . In (3.6), the matrix product  $B_{nm}(q)B_{ni}(p^2z_\rho)B_{nm}(q)$  was written as the band Toeplitz matrix  $B_{n,i+2m}(q^2p^2z_\rho)$ , generated by the function  $q^2p^2z_\rho$ , plus the low rank correction matrix  $\Delta$ .

It is known [2] that the matrix  $\Delta$  is given by

$$\Delta = B_{nm}(q)H(q)H(p^2z_\rho) + B_{nm}(q)H^R(q)H^R(p^2z_\rho) + H(q)H(qp^2z_\rho) + H^R(q)H^R(qp^2z_\rho),$$

where  $H(q)$ ,  $H(p^2z_\rho)$  and  $H(qp^2z_\rho)$  are Hankel matrices produced by the trigonometric polynomials  $q$ ,  $p^2z_\rho$  and  $qp^2z_\rho$ , respectively, while  $H^R$  denotes the matrix obtained from  $H$  by reversing the order of its rows and columns.

It is obvious that  $\Delta$  is a low rank correction matrix that has nonzero elements



only in the upper left and lower right triangles as this is illustrated below

$$\Delta = \begin{pmatrix} * & \cdots & * & 0 & \cdots & 0 \\ \vdots & \ddots & 0 & \ddots & 0 & \vdots \\ * & 0 & \ddots & 0 & & 0 \\ 0 & \ddots & 0 & \ddots & 0 & * \\ \vdots & 0 & & 0 & \ddots & \vdots \\ 0 & \cdots & 0 & * & \cdots & * \end{pmatrix}.$$

It is clear that the elements of  $\Delta$  are bounded and the size of the triangles depends only on the bandwidths  $m$  and  $\hat{l}$  and are independent of  $n$ .

It remains to prove that  $\rho(B_{n\hat{l}}^{-1}(p^2 z_\rho)\Delta)$  is bounded. For this, we write the matrices in the following block forms

$$B_{n\hat{l}}^{-1}(p^2 z_\rho) = \begin{pmatrix} B_1 & * & B_2 \\ * & * & * \\ B_2^T & * & B_1^R \end{pmatrix}, \quad \Delta = \begin{pmatrix} D & & \\ & O & \\ & & D^R \end{pmatrix},$$

where  $B_1, B_2$  are  $k \times k$  matrices if  $D$  has  $k$  nonzero anti-diagonals.

Since the only nonzero columns of the matrix  $B_{n\hat{l}}^{-1}(p^2 z_\rho)\Delta$  are its first  $k$  and last  $k$  ones, the nonidentically zero eigenvalues of  $B_{n\hat{l}}^{-1}(p^2 z_\rho)\Delta$  will be the eigenvalues of the matrix

$$\begin{pmatrix} B_1 D & B_2 D^R \\ B_2^T D & B_1^R D^R \end{pmatrix}.$$

In view of Lemma 3.2 this matrix is bounded and so are its eigenvalues which proves the present statement.  $\square$

So, the eigenvalues that are greater than  $h_{\max}$ , have an upper bound. An open question remains regarding the eigenvalues that may lie in the interval  $(0, h_{\min})$ . However, strong numerical evidence suggests that in the spectrum of the preconditioned matrix obtained by our approach (see Figures 5.1, 5.2, 5.3), these eigenvalues have a lower bound independent of  $n$ . Moreover, as one can see from Figures (5.1(b)-(d), 5.2(b), 5.3(b)), the out of the main interval eigenvalues appear in pairs. In addition, the elements of each pair tend to each other as  $n$  tends to infinity. In view of this observation the convergence analysis of the PCG method in [1] assures us that our method will not be seriously affected and the convergence of it will remain superlinear which is the optimal cost for this method.

**4. Computational analysis and modifications of the method.** In this section we will try to compare, from the computational point of view, our preconditioner with the most recent band-Toeplitz preconditioner proposed in [15]. The latter has in general the best performance from all the previous ones, when the generating function  $f$  is nonnegative and has zeros of even order.

The main computational cost in every PCG iteration is due to the Toeplitz matrix-vector product  $T_n(f)x$  and to the solution of a system with coefficient matrix the preconditioner itself. The first one is the same for both methods and can be computed by means of Fast Fourier Transform (FFT) in  $30(n \log 2n)$  operations (ops) in



a sequential machine or in  $O(\log 2n)$  steps in the parallel PRAM model of computation, when  $O(n)$  processors are used. For the inversion of the preconditioners things slightly change. If we use band Toeplitz preconditioners then their halfbandwidth  $\hat{l}_1$  represents the degree  $l_1$  of the Chebyshev approximation plus the degree  $\rho$  of the trigonometric polynomial which eliminates the zeros of  $f$ . The inversion of such type of matrices can be achieved using the  $LDL^T$  factorization method in  $n(\hat{l}_1^2 + 8\hat{l}_1 + 1)$  ops. We mention that this method is preferable from the band Cholesky procedure because the latter requires the computation of  $n$  square roots, which is quite expensive when  $n$  is large.

In the case of our preconditioner the inversion requires two band matrix vector products of total cost  $n(4m + 2)$  ops, where  $m$  is the halfbandwidth and coincides with the degree of the denominator in the rational approximation. In addition, the inversion of  $B_{n,\hat{l}_2}$ , as in the previous case, can be performed in  $n(\hat{l}_2^2 + 8\hat{l}_2 + 1)$  ops, where  $\hat{l}_2 = \rho + 2l_2$  and  $l_2$  represents the degree of the numerator of the rational approximation. So the total cost per iteration for this step of the algorithm of the PCG method is about

$$Cost_{it} = n(\hat{l}_2^2 + 8\hat{l}_2 + 4m + 3).$$

When  $n$  is large, the complexity of the method is strongly dominated by the first step which requires  $O(n \log 2n)$  ops and the methods are essentially equivalent in complexity per iteration. Thus the costs of finding  $B_{n,\hat{l}_1}^{-1}$  and  $B_{n,m} B_{n,\hat{l}_2}^{-1} B_{n,m}$ , where  $l_1 = l_2 + m$ , are comparable.

In case  $n$  is not large enough, taking  $l_2 = \frac{l_1}{2} - 1$  and making some calculations, we can see that the two preconditioning strategies are approximately equivalent even when  $m = \rho l_1$ .

According to this observation, if we have two candidates of rational approximations of  $f$  with almost the same relative error and degrees  $(l_1, m_1)$ ,  $(l_2, m_2)$  with  $l_1 + m_1 \approx l_2 + m_2$ , it is preferable, from the computation point of view, to choose as the generating function for our preconditioner the one which has the larger  $m$  and smaller  $l$ .

Finally, we will focus on the calculation of rational approximation of degree  $(l, m)$  of a positive continuous function  $f$ . In the recent literature many different strategies that produce this kind of approximation [11] can be found. Each of them is most suitable for certain classes of functions but the one which is based on the Remez algorithm seems to be, in general, quite efficient for a large variety of functions. The starting point of this category of algorithms is to construct a rational approximation using rational interpolation and then this rational approximation is used to generate a better approximation until an alternative set of  $m + l + 2$  points is achieved. This procedure consists of adjusting the choice of the interpolation points in such a way as to ensure that the relative error decreases. In practice this method can fail in some cases. Usually, problems are caused either from the fact that the extreme values of the relative error occur more than  $m + l + 2$  times, or the starting rational interpolation has zeros in the interval in which this approximation is sought. The first difficulty is usually overcome by seeking a rational approximation of a different degree or by designing a more robust algorithm. A trick that often works in the latter case is, instead of asking again for a rational approximation of a different degree, to start with an approximation that is valid over a shorter interval and use it as a starting point for an approximation on a slightly larger interval. Iterative application of this procedure may enable us to obtain a final approximation in the desired interval.

TABLE 5.1  
Number of iterations for  $f_1(x)$

n	$B_n^{*1}$	$\hat{B}_n^1$	$B_n^{*3}$	$\hat{B}_n^3$	$B_n^{*4}$	$\hat{B}_n^4$	$M_n^{0,1}$	$R_n^{0,1}$	$M_n^{1,1}$	$R_n^{1,1}$	$M_n^{1,2}$	$R_n^{1,2}$
16	9	8	9	7	7	6	8	7	6	6	5	5
32	10	10	11	8	9	7	10	9	7	7	6	6
64	13	12	11	10	9	8	11	11	9	9	8	8
128	15	15	12	11	10	10	12	13	11	11	10	10
256	16	16	12	13	10	10	13	13	12	12	11	11
512	16	16	13	13	10	11	13	14	13	13	11	12

For the convergence rate of the approximation method we can not give a theoretical result, but the facts that its computational cost is independent of  $n$  and the computations are done only once for a given function make us believe that this problem does not play an important role in the whole procedure.

**4.1. Modifications of the method.** The idea of constructing a preconditioner from a rational approximation of a function can be used in exactly the same way in case of rational interpolation at the Chebyshev points. The advantage of this modification is the easiness of its calculation. Nevertheless, it is worth noticing that we can not assure that this interpolation would not have zeros in the interval of approximation. Despite this, whenever the preconditioning gives us poor results, this technique may give, at least for certain classes of  $f$ , results similar to the corresponding ones by the best Chebyshev approximation.

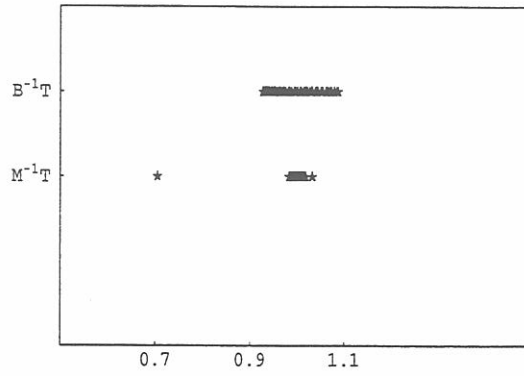
Another modification of this kind of preconditioning would be the following. First, we approximate the function  $\frac{f}{z_\rho}$  by a rational approximation  $\frac{p_l}{r_k}$ , where  $k$  can be very large. Then we approximate the function  $\sqrt{r_k}$  using a polynomial Chebyshev approximation  $q_m$ . Finally, the ratio  $\frac{p_l}{q_m}$  is considered as an approximation of  $\frac{f}{z_\rho}$ . So, the preconditioner matrix  $\tilde{M}$  for the solution of (1.1) would be

$$(4.1) \quad \tilde{M}_n = B_{nm}^{-1}(q) B_{ni}(pz) B_{nm}^{-1}(q), \quad \hat{l} = l + \rho,$$

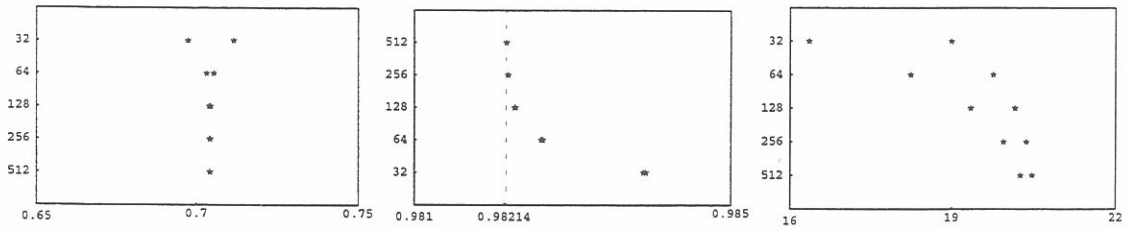
instead of  $M_n$  in (3.1). After this, all the previous theory developed holds the same.

The main point of this method is to approximate directly  $\frac{f}{z_\rho}$  instead of  $\sqrt{\frac{f}{z_\rho}}$  and possibly with a polynomial of higher degree in the denominator. Then considering that this can take care of every possible abnormalities of  $f$ , we approximate the denominator by a polynomial of lower degree by the Chebyshev technique. We remark here, that numerical experiments show that this matrix is not in general so good as a preconditioner compared with  $M_n$  or with the band-Toeplitz preconditioner obtained in [15]. This is because we make approximations in two levels. First, we take the rational approximation and then the Chebyshev approximation of the square root of the denominator of the first approximation. So, the overall approximation error seems to become much larger.

**5. Numerical examples and concluding remarks.** In this section, we present some numerical examples. The aim of these examples is twofold: i) to show, by numerical evidence, the correctness of our observations regarding the asymptotical spectral analysis of the preconditioned matrices and ii) to compare the convergence rate of our preconditioner with that of the band Toeplitz preconditioner proposed in [15]. We



(a) The main mass of the eigenvalues of the preconditioned matrices



(b) The lower extreme pair

(c) The second upper pair

(d) The upper extreme pair

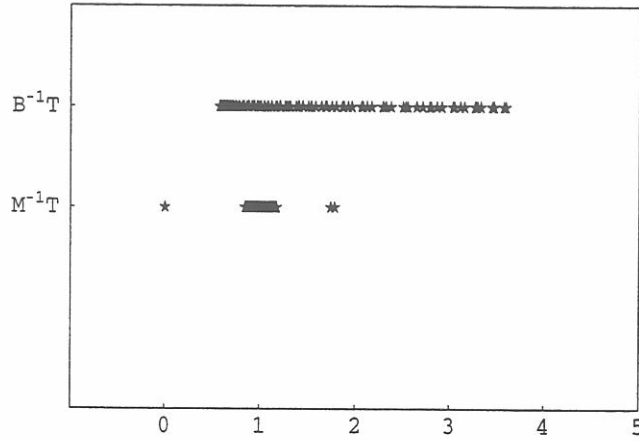
FIG. 5.1. Spectra of  $(M_n^{2,2})^{-1}T_n(f_1)$  and  $(B_n^{*5})^{-1}T_n(f_1)$  for  $n = 128$  and behavior of the pairs of eigenvalues that lie outside the interval  $[h_{\min}, h_{\max}]$  with  $h_{\min} = 0.98214$

TABLE 5.2  
Number of iterations for  $f_2(x)$

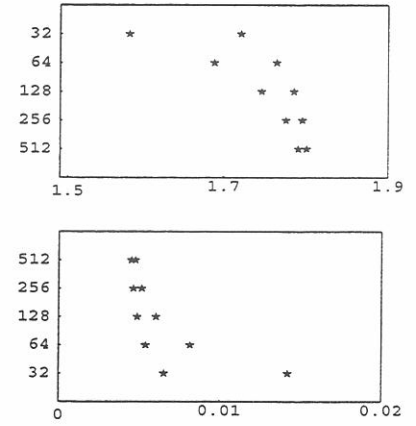
n	$B_n^{*3}$	$B_n^{*4}$	$B_n^{*5}$	$B_n^{*6}$	$M_n^{1,1}$	$R_n^{2,2}$
16	8	8	7	8	8	6
32	13	13	12	11	11	7
64	19	18	15	13	12	9
128	24	19	17	14	12	11
256	25	21	18	15	13	13
512	27	22	18	16	14	14

use the latter to compare it with ours because it is the most efficient technique for preconditioning Toeplitz matrices generating by functions with zeros of even order. Our test functions are the following

$$i) f_1(x) = x^4, \quad ii) f_2(x) = \frac{2x^4}{1 + 25x^2}$$

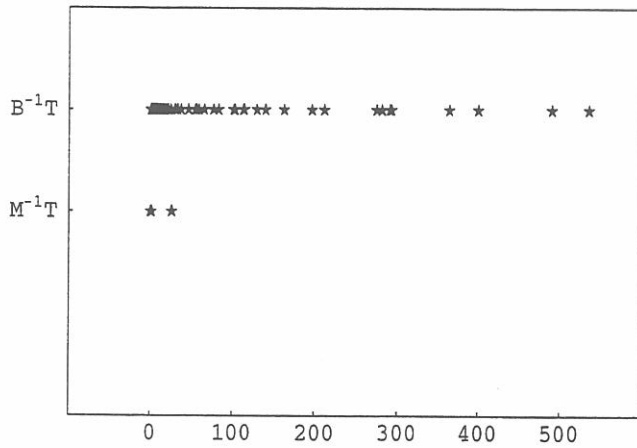


(a)

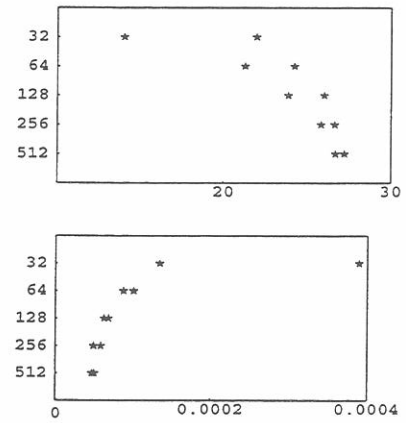


(b) The two pairs of extreme eigenvalues

FIG. 5.2. Spectra of  $(M_n^{1,1})^{-1}T_n(f_2)$  and  $(B_n^{*3})^{-1}T_n(f_2)$  for  $n = 128$  and behavior of the pairs of eigenvalues that lie outside the interval  $[h_{min}, h_{max}]$



(a)



(b) The two pairs of extreme eigenvalues

FIG. 5.3. Spectra of  $(M_n^{1,2})^{-1}T_n(f_3)$  and  $(B_n^{*3})^{-1}T_n(f_3)$  for  $n = 256$  and behavior of the pairs of eigenvalues that lie outside the interval  $[h_{min}, h_{max}]$

and

$$iii) \quad f_3(x) = \begin{cases} (x-3)^4(x-1)^2 & 0 \leq x \leq \pi, \\ (x+3)^4(x+1)^2 & -\pi \leq x \leq 0. \end{cases}$$

An effort was made to choose functions of different behaviors which produce ill-

TABLE 5.3  
Number of iterations for  $f_3(x)$ .

n	$B_n^{*3}$	$B_n^{*5}$	$B_n^{*7}$	$M_n^{1,2}$	$R_n^{(1,2)}$
16	9	7	7	9	8
32	17	14	13	18	11
64	34	28	22	21	14
128	65	48	36	21	20
256	111	69	54	23	24
512	152	93	66	23	27

conditioned matrices  $T_n$ . The Toeplitz matrices produced have Euclidean condition numbers of order  $O(n^4)$ . In our experiments we solve the system  $T_n(f)x = b$  where  $b$  is the vector having all its components equal to one. As a starting initial guess of solution the zero vector is used and as a stopping criterion the validity of  $\frac{\|r_k\|_2}{\|r_0\|_2} \leq 10^{-7}$  is considered, where  $r_k$  is the residual vector after  $k$  iterations. The matrices and the rational approximations were performed using Mathematica in order to have more accurate results while all the other computations were performed using Matlab.

In the Tables we report the number of iterations needed until convergence is achieved in each case,  $B_n^{*l}$  denotes the optimal band Toeplitz preconditioner [15] which is generated by the trigonometric polynomial  $z_\rho g_l$ , with  $g_l$  being the best Chebyshev approximation of  $\frac{f}{z_\rho}$  out of  $\mathcal{P}_l$ ,  $\hat{B}_n^l$  is the band Toeplitz preconditioner where  $\hat{g}_l$  is the interpolation polynomial at the Chebyshev points,  $M_n^{l,m}$  denotes our main proposed preconditioner obtained by the best rational approximation procedure of degree  $(l, m)$  and  $R_n^{l,m}$  denotes the preconditioner that results after applying rational interpolation of degree  $(l, m)$ .

In Figures 5.1(a), 5.2(a), 5.3(a), the spectra of the matrices  $M_n^{-1}T_n(f_i)$ ,  $i = 1, 2, 3$ , are illustrated, while in 5.1(b)-(d), 5.2(b), 5.3(b) we focus on the behavior of the pairs of eigenvalues of the matrix lying outside the interval  $[h_{\min}, h_{\max}]$  for different values of  $n$ . The boundness and the convergence in pairs is obvious in all figures. Especially, we stress the case of figures (5.1) and (5.3) where as we expected from the theory at most eight eigenvalues would lie outside the interval  $[h_{\min}, h_{\max}]$  but in practice, for the first test function, only three pairs of eigenvalues lie outside this interval, one of which (the second lower pair) moves very close to the lower bound  $h_{\min} = 0.98214$  while, for the third test function, only two pairs lie outside this interval. Finally, we remark that in the case of  $f_3$  and for  $n = 512$ , the preconditioning by band Toeplitz  $B^{*3}$  "clusters" the eigenvalues of the preconditioned matrix in  $[0.5, 584.3]$ ,  $B^{*5}$  in  $[0.36, 104.7]$  while  $M^{1,2}$  collects the main mass of them in  $[0.67, 1.65]$  and  $R^{1,2}$  collects it in  $[0.95, 14.25]$ .

## REFERENCES

- [1] O.AXELSSON AND G.LINDSKÖG, *On the rate of convergence of the preconditioned conjugate gradient method*, Numer. Math., 52 (1986), pp. 499-523.
- [2] A. BÖTTCHER AND B. SILBERMANN, *Introduction to Large Truncated Toeplitz Matrices*, Springer Verlag, 1998.

- [3] J. R. BUNCH, *Stability of methods for solving Toeplitz systems of equations*, SIAM J. Sci. Stat. Comput., 6 (1985), pp. 349-364.
- [4] R. CHAN, *Toeplitz Preconditioners for Toeplitz Systems with Nonnegative Generating Functions*, SIAM J. of Numer. Anal., 11 (1991), pp. 333-345.
- [5] R. CHAN AND P. TANG, *Fast Band-Toeplitz preconditioners for Hermitian Toeplitz systems*, SIAM J. Sci. Comp., 15 (1994), pp. 164-171.
- [6] F. DI BENEDETTO, *Analysis of Preconditioning Techniques for ill-conditioned Toeplitz matrices*, SIAM J. Sci. Comput., 16 (1995), pp. 682-697.
- [7] F. DI BENEDETTO AND S. SERRA, *A unifying approach to abstract matrix algebra preconditioning*, Numer. Math., 82 (1999), pp. 57-90.
- [8] U. GRENANDER AND G. SZEGÖ, *Toeplitz Forms and Their Applications*, 2nd edition, Chelsea, New York, 1984.
- [9] X.A. JIN, *Hartley Preconditioners for Toeplitz Systems generated by positive continuous functions*, BIT, 34 (1994), pp. 367-371.
- [10] G. LORENTZ, *Approximation of Functions*, 2nd edition, Chelsea, New York, 1986.
- [11] M.J.D. POWELL, *Approximation theory and methods*, Cambridge Univ. Press, 1982.
- [12] T. RIVLIN, *Introduction to the Approximation of functions*, Dover Pubs, 1981.
- [13] S. SERRA, *New PCG based algorithms for the solution of Hermitian Toeplitz systems*, Calcolo, 32 (1995), pp. 154-176.
- [14] S. SERRA, *On the extreme spectral properties of Toeplitz matrices generated by  $L^1$  functions with several minima (maxima)*, BIT, 36 (1996), pp. 135-142.
- [15] S. SERRA, *Optimal, Quasi-Optimal and Superlinear Band-Toeplitz preconditioners for asymptotically ill-conditioned positive definite Toeplitz Systems*, Math. Comp., 66 (1997), pp. 651-665.
- [16] S. SERRA, *Toeplitz preconditioners constructed from linear approximation processes*, SIAM J. Matrix. Anal. Appl., 20 (1998), pp. 446-465.
- [17] S. SERRA, *A Korovkin-type Theory for finite Toeplitz operators via matrix algebras*, Numer. Math., 82 (1999), pp. 117-142.
- [18] G. STRANG, *A proposal for Toeplitz matrix calculation*, Stud. Appl. Math., 74 (1986), pp. 171-176.
- [19] H. WIDOM, *Toeplitz Matrices*, In Studies in real and Complex analysis, I. Hirshman Jr. Ed., Math. Ass. Am., 1965.
- [20] J.H. WILKINSON, *The Algebraic Eigenvalue Problem*, Oxford Press, Oxford, 1965.