

OPTIMAL, QUASI-OPTIMAL AND SUPERLINEAR  
BAND-TOEPLITZ PRECONDITIONERS FOR  
ASYMPTOTICALLY ILL-CONDITIONED POSITIVE DEFINITE  
TOEPLITZ SYSTEMS

STEFANO SERRA

ABSTRACT. In this paper we are concerned with the solution of  $n \times n$  Hermitian Toeplitz systems with nonnegative generating functions  $f$ . The preconditioned conjugate gradient (PCG) method with the well-known circulant preconditioners fails in the case where  $f$  has zeros. In this paper we consider as preconditioners band-Toeplitz matrices generated by trigonometric polynomials  $g$  of fixed degree  $l$ . We use different strategies of approximation of  $f$  to devise a polynomial  $g$  which has some analytical properties of  $f$ , is easily computable and is such that the corresponding preconditioned system has a condition number bounded by a constant independent of  $n$ . For each strategy we analyze the cost per iteration and the number of iterations required for the convergence within a preassigned accuracy. We obtain different estimates of  $l$  for which the total cost of the proposed PCG methods is optimal and the related rates of convergence are superlinear. Finally, for the most economical strategy, we perform various numerical experiments which fully confirm the effectiveness of approximation theory tools in the solution of this kind of linear algebra problems.

1. INTRODUCTION

The aim of this paper is to introduce and analyze new strategies for the solution by PCG method of  $n \times n$  Hermitian Toeplitz systems [20, 21]

$$A_n \mathbf{x} = \mathbf{b}.$$

Toeplitz matrices are assumed to be generated by  $2\pi$ -periodic integrable real-valued functions  $f$  defined on the fundamental interval  $[-\pi, \pi]$ , in the sense that the coefficients of  $A_n$  are given by the Fourier coefficients  $a_m$  of  $f$ : more precisely we have

$$[A_n]_{j,k} = a_{j-k} = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) e^{-i(j-k)x} dx, \quad 0 \leq j, k \leq n-1.$$

We point out that the generating function  $f$  is given in some applications of Toeplitz systems. Classical examples are the kernels of the Wiener-Hopf equations [18], the spectral density functions in stationary stochastic processes [20] and the point-spread functions in image processing [24].

---

Received by the editor January 20, 1995 and, in revised form, January 26, 1996.

1991 *Mathematics Subject Classification*. Primary 65F10, 65F15.

*Key words and phrases*. Toeplitz matrix, Chebyshev interpolation and approximation, conjugate gradient method, Remez algorithm.

If the generating function is continuous and positive there are many types of preconditioners [10, 15, 16, 6] such as circulant matrices [14],  $\tau$  matrices [5], Hartley matrices [6]; these preconditioners lead to superlinearly convergent PCG methods.

When  $f$  has zeros, i.e.,  $\text{ess inf } f = 0$  we know [20] that the Euclidean condition number  $\mu_2(A_n(f))$  of  $A_n = A_n(f)$  grows to  $\infty$  for  $n$  tending to  $\infty$ ; in [30, 31] estimates of  $\mu_2(A_n(f))$  as a function of  $n$  and of the order  $r$  of the “zeros” of  $f$  are given. In the case where  $r = 2k$  is an even number only  $\tau$  preconditioners [15] and band-Toeplitz preconditioners [7, 16] are shown to be able to reduce the condition number from  $O(n^{2k})$  to  $O(1)$ . More general statements and strategies useful to handle the case where  $f$  has also zeros of odd or noninteger orders can be found in [28, 32].

The main idea (see [16]) is to find a trigonometric polynomial  $g$  for which  $r < f/g < R$  where  $r, R$  are positive constants. The associated band-Toeplitz matrix  $A_n(g)$  results to be the desired preconditioner in the sense that the spectrum of  $A_n^{-1}(g)A_n(f)$  lies in  $(r, R)$  for any dimension  $n$ .

The quoted idea resulted to be very flexible and, actually, has been successfully applied to the case of nondefinite Toeplitz problems [29], block Toeplitz problems [27] and, joint with circulant structures, to the case of non-Hermitian Toeplitz problems [8].

More recently, R. Chan and P. Tang [11] have proposed to increase the bandwidth of  $A_n(g)$  to get extra degrees of freedom. They calculate  $g$  by means of the Remez algorithm by minimizing  $h = \|(f-g)/f\|_\infty$  over all the polynomials  $g$  of fixed degree  $l$ . In this paper we prove that this minimization property enables one not only to match the zeros of  $f$  but also to minimize  $R/r$  obtaining (by using Theorem 3.1 in [16]) the best band-Toeplitz preconditioner in the class of all the band-Toeplitz matrices of fixed bandwidth  $2l + 1$ . Moreover we perform a more accurate analysis than [11] of the convergence properties of the preconditioned systems defined in [11].

Since the Remez algorithm can be heavy from a computational point of view, we propose two new techniques to minimize “in a certain sense”  $(f-g)/f$ . These strategies are such that  $g$  is easier to calculate (for example, for one of the proposed polynomials  $g = g^B$  we use only few Fast Fourier Transforms (FFT) of order  $l-k$ ) and the preconditioned systems have an  $O(1)$  condition number for which we can exhibit upper bounds depending on  $l, n$  and on the “regularity” features of  $f$ .

Therefore, we can estimate the number of iterations to reach the solution within a preassigned accuracy  $\epsilon$ ; on the other hand, the solution of a system  $A_n(g)\mathbf{y} = \mathbf{c}$  can be obtained in  $O(l^2n)$  arithmetic operations (ops), by using a classic band solver [19], or in  $O(ln)$  ops [17] (see also [9]).

Hence, balancing the cost of a single iteration of the PCG and the number of required iterations, it is possible to estimate the optimal bandwidth  $l$ , which allows to minimize the *total amount* of calculations to reach the solution of  $A_n(f)\mathbf{x} = \mathbf{b}$  within a preassigned tolerance  $\epsilon$ .

The outline of the paper is the following. In section 2 we analyze the convergence rate of the PCG method proposed in [11]. In sections 3 and 4 we introduce two new preconditioners and we perform a study of the convergence properties of our PCG methods. In the subsequent section 5, first we discuss the cost of the different PCG methods and then we indicate how to estimate the value  $l_{opt}$  such that the global optimal preconditioner has to be searched in the band-Toeplitz matrices of bandwidth  $2l_{opt} + 1$ . In section 6 we observe that, by choosing  $l$  as special functions

of  $n$ , we may construct superlinearly convergent PCG methods having a total cost of  $O(n \log n)$  ops. Finally in the last section we perform several numerical experiments showing the effectiveness of the proposed ideas.

2. CONVERGENCE ANALYSIS OF THE PCG METHOD  
OF R. CHAN AND P. TANG

We study the convergence speed of the PCG proposed in [11] in terms of the generating functions  $f$  and  $g$ .

Firstly we recall some known results.

**Theorem 2.1.** *Let  $m_f$  and  $M_f$  be the essinf and the esssup of  $f$  in  $[-\pi, \pi]$ . If  $m_f < M_f$  then  $\forall n > 0$  we have*

$$m_f < \lambda_i(A_n(f)) < M_f$$

where  $\lambda_i(X)$  is the  $i$ -th eigenvalue of  $X$  arranged in nondecreasing order. If  $m_f \geq 0$  then  $A_n(f)$  is positive definite.

*Proof.* See [20, 7]. □

**Theorem 2.2.** *Let  $f, g \in L^1[-\pi, \pi]$  be functions essentially nonnegative, i.e.,  $m_f, m_g \geq 0$ . The matrices  $A_n(f), A_n(g)$  are positive definite (see Theorem (2.1)) and the eigenvalues  $\lambda_i^n$  of  $A_n^{-1}(g)A_n(f)$  arranged in nondecreasing order are such that:*

1.  $\lambda_i^n \in (r, R)$ ,  $r, R$  being the ess inf and the ess sup of  $f/g$ , respectively.
2.  $\bigcup_{n \in \mathbb{N}} \bigcup_{i \leq n} \lambda_i^n$  is dense in the “essential range”  $\mathcal{ER}(f/g)$  of  $f/g$  (the essential range of an integrable function  $h$  defined on  $I$  is the set of all  $y$  real numbers for which,  $\forall \epsilon > 0$  the Lebesgue measure of  $\{x \in I : h(x) \in (y - \epsilon, y + \epsilon)\}$  is positive [28]).
3.  $\lim_{n \rightarrow \infty} \lambda_1^n = r, \quad \lim_{n \rightarrow \infty} \lambda_n^n = R.$

*Proof.* Under the assumption that  $f, g$  and  $f/g$  are continuous the claimed thesis follows from Theorems 3.1 and 3.2 in [16]. Note that in the case where  $f/g$  is continuous the essential range of  $f/g$  coincides with  $[r, R]$  and, in general, if  $f/g$  is piecewise continuous the set  $\mathcal{ER}(f/g)$  is the closure of the usual image of  $f/g$ .

Under the weaker hypothesis that  $f$  and  $g$  are only integrable follows from Theorem 2.2 in [30] and Theorem 3.1 in [28]. □

*Remark 1.* The third statement of the preceding theorem has also been proved in [12]. However, we notice that Theorem 2.2 is much more powerful, since it indicates a “global property” of distribution of the eigenvalues. For instance, a consequence of the second part of the considered theorem is that for any nonnegative integer  $k$  fixed with respect to the dimension  $n$  we find the following limit relations

$$\lim_{n \rightarrow \infty} \lambda_k^n = r, \quad \lim_{n \rightarrow \infty} \lambda_{n-k}^n = R.$$

In addition, we can conclude that the spectrum of the preconditioned matrix  $A_n^{-1}(g)A(f)$  is, for large  $n$ , “uniformly distributed” in the image of  $\frac{f}{g}$ . This means

[32] that the set  $\left\{ \frac{f}{g} \left( \frac{2\pi j}{n} \right) \right\}_{j=1}^n$ , suitably ordered, describes asymptotically the set  $\{\lambda_j^n\}_{j=1}^n$ .

Finally we stress that these results are useful in order to understand very precisely the convergence rates of PCG methods based on Toeplitz preconditioners.

Actually, owing to the sophisticated results [3] about the convergence speed of PCG algorithms, we conclude that the knowledge of the asymptotical behaviour of  $\lambda_1^n$  and  $\lambda_n^n$  is not only useful (part 3 or [12]), but also the global distribution of the preconditioned spectrum [16, 29, 31].

By means of Theorem 2.1 in [11] it is shown that, if  $g$  is a polynomial of degree  $l$  such that

$$(1) \quad \left\| \frac{f-g}{f} \right\|_{\infty} = h < 1,$$

then  $A_n(g)$  is positive definite and the Euclidean condition number of  $A_n^{-1/2}(g) \cdot A_n(f) A_n^{-1/2}(g)$  is bounded by a positive constant, i.e.,

$$\mu_2 \left( A_n^{-1/2}(g) A_n(f) A_n^{-1/2}(g) \right) \leq \frac{1+h}{1-h}.$$

Consequently, by standard error analysis of the PCG method [2], Chan and Tang conclude that the number of iterations for convergence within a tolerance  $\epsilon$  is bounded by

$$N(h, \epsilon) = \frac{1+h}{2(1-h)} \log \left( \frac{1}{\epsilon} \right) + 1.$$

In the following theorem we refine the result in [11] and we show that the former bound is a very sharp bound, i.e., the number of iterations that we expect cannot be much less than  $N(h, \epsilon)$ .

**Theorem 2.3.** *Let  $f \geq 0$  be a continuous function and  $g$  be a polynomial of degree  $l$  such that the relative error  $h$  is less than 1. Then*

1.  $\mu_2^n = \mu_2 \left( A_n^{-1/2}(g) A_n(f) A_n^{-1/2}(g) \right) < \frac{1+h}{1-h}$ .
2.  $\bigcup_{n \in \mathbb{N}} \bigcup_{i \leq n} \lambda_i^n$  is dense in  $[1/(1+h), 1/(1-h)]$ ,  $\lambda_i^n$  being the  $i$ -th eigenvalue of  $A_n^{-1}(g) A_n(f)$ . Therefore we have no clusters in  $(1/(1+h), 1/(1-h))$  but practically a “uniform distribution” of the spectrum of the preconditioned matrix.
3.  $\lim_{n \rightarrow \infty} \mu_2^n = (1+h)/(1-h)$ .

*Proof.* From  $\left\| \frac{f-g}{f} \right\|_{\infty} = h < 1$  we deduce that

$$-h \leq 1 - \frac{g}{f} \leq h, \quad g \geq 0,$$

and, consequently,

$$(2) \quad \frac{1}{1+h} \leq \frac{f}{g} \leq \frac{1}{1-h}.$$

Now we may apply Theorem 2.2 obtaining (1), (2) and (3). Moreover, it is worth pointing out that (2) implies that  $f$  has only zeros of even order, because of the fact that  $g$  is a nonnegative trigonometric polynomial.  $\square$

As a final remark of this section we can state the following property.

**Theorem 2.4.** *The preconditioner  $A_n(g^*)$ , where  $g^*$  is the best relative Chebyshev approximation of  $f$  of degree  $l$ , is optimal in the sense that  $N(h, \epsilon)$  is minimal for  $g = g^*$ .*

*Proof.*

$$h^* = \left\| \frac{f - g^*}{f} \right\|_\infty = \min_{g \in \mathbf{P}_l} \left\| \frac{f - g}{f} \right\|_\infty,$$

$\mathbf{P}_l$  being the class of the trigonometric polynomials of degree at most  $l$ . As for  $h \in (0, 1)$  the function  $N(h, \epsilon)$  is an increasing function of  $h$ , it is trivial to note that the minimal value of  $N(h, \epsilon)$  is attained for  $h = h^*$ .

Observe that we can suppose  $h^* < 1$  as proved in the third part of Theorem 4.1 in section 4. □

### 3. NEW PRECONDITIONING STRATEGIES

We start this section with the following observation: in the case where  $f$  is non-negative and has some zeros in  $[-\pi, \pi]$ , band-Toeplitz preconditioners can reduce the condition number to a value uniformly bounded by a constant independent of the dimension  $n$  only when the zeros of  $f$  have even order [7, 16].

Actually, since a nonnegative trigonometric polynomial  $g$  can have only zeros of even order, in light of Theorem 2.2 it is trivial to conclude that the union of the spectra of  $A_n^{-1}(g)A_n(f)$  cannot be contained in a positive interval in the case where  $f$  has a zero of order  $r \neq 2q$ , for any positive integer  $q$ . Therefore, in the following, we assume that  $f$  is continuous and has only zeros of even order.

Now we define by  $z_k$  the polynomial of minimum degree  $k$  containing all the zeros of  $f$  with their orders and the generating function  $g$  of our preconditioners in the following way:

$$g = z_k g_{l-k}, \quad \text{degree}(g) = l \geq k.$$

$g_{l-k}$  is a trigonometric polynomial of degree  $l - k$  and can be chosen, for example, in light of these two strategies.

**A:**  $g_{l-k}$  is the best Chebyshev approximation of  $\hat{f} = f/z_k$ , i.e.,

$$\|\hat{f} - g_{l-k}\|_\infty = \min_{g \in \mathbf{P}_{l-k}} \|\hat{f} - g\|_\infty.$$

**B:**  $g_{l-k}$  is the trigonometric polynomial of degree at most  $l - k$  interpolating  $\hat{f}$  at the  $l - k + 1$  zeros of the  $(l - k + 1)$ -th Chebyshev polynomial of the first kind.

Observe that we cannot choose  $g$  directly as the best Chebyshev approximation of  $f$  for two reasons: we are not guaranteed that  $g$  is nonnegative since  $f$  has zeros. In fact, if  $g = g_l$  such that  $\|f - g_l\|_\infty = \min_{g \in \mathbf{P}_l} \|f - g\|_\infty = E_l(f)$ ,  $l = \text{degree}(g_l)$ , then we have

$$(3) \quad f(x) - E_l(f) \leq g(x) \leq f(x) + E_l(f)$$

and  $\lim_{l \rightarrow \infty} E_l(f) = 0$ , but, since  $f$  has zeros, it may happen that  $g$  assumes negative values (see equation (3)) in a suitable neighbourhood of each zero of  $f$ . Consequently, by virtue of the classical spectral theory on Toeplitz matrices,  $A_n(g)$  is not positive definite for any  $n$  large enough and cannot be used as preconditioner [20]. In addition, it is worth pointing out that relation (3) does not imply that  $f/g$  and  $g/f$  are bounded because, in general,  $g$  has different zeros with respect to  $f$ . Finally, from Theorem 2.2, we cannot expect a convergence speed independent of the dimension  $n$ , since  $r = 0$  and/or  $R = +\infty$ .

From a computational point of view we remark that  $g_{l-k}$  in **(A)** can be calculated by using the standard Remez algorithm [26] with respect to the classical trigonometric basis  $\{1, \sin(qx), \cos(qx)\}$ , while the calculation of  $g^*$  in [11] is performed by using a modified version of the Remez algorithm [33] with the basis  $\{1/f(x), \sin(qx)/f(x), \cos(qx)/f(x)\}$ ; in this case it is possible to observe instability problems due to the fact that  $f$  has zeros (see section 6).

For the calculation of  $g_{l-k}$  in **(B)**, on the other hand, we have no problems: this polynomial can be calculated, very easily, with few FFTs of order  $l - k$  by means of a classical trigonometric representation of the interpolating polynomial at Chebyshev zeros (see section 6).

#### 4. CONVERGENCE ANALYSIS

We perform a convergence analysis of the PCG methods proposed in the former sections; it is worth pointing out that this analysis gives further information on the convergence properties of the PCG method proposed by R. Chan and P. Tang.

First we introduce a result which makes a link between Theorem 2.2 in [11] and Theorem 2.2, i.e., Theorems 3.1 and 3.2 in [16].

**Theorem 4.1.** *Let  $f$  be a nonnegative continuous function defined in  $[-\pi, \pi]$  with zeros of even order, then the following statements hold.*

1. *There exists a nonnegative trigonometric polynomial  $z_k$  of minimal degree  $k$  [7, 16] such that*

$$0 < r_k < \frac{f}{g} < R_k < \infty.$$

2. *If  $g$  is a trigonometric polynomial such that  $f/g \in (r, R)$  (one of the hypotheses of Theorem 2.2),  $r, R$  being positive constants, then there exists  $\alpha > 0$  for which*

$$h_{\alpha g} = \left\| \frac{f - \alpha g}{f} \right\|_{\infty} < 1.$$

3. *If  $l \geq k$  (with  $k$  the degree of the minimal polynomial  $z_k$ ), then*

$$h^* = \min_{g \in \mathcal{P}_l} \left\| \frac{f - g}{f} \right\|_{\infty} < 1.$$

(Recall that  $h^* < 1$  is one among the hypotheses of the main theorem in [11].)

*Proof.* Let  $x_1, \dots, x_j$  be the zeros of  $f$  and  $2l_1, \dots, 2l_j$  be the orders of such zeros. The linear polynomial  $2 - 2\cos(x - \hat{x})$  is nonnegative and is clearly the polynomial of minimal degree which has in  $\hat{x}$  a zero of order 2. Consequently  $z_k$  is easily constructed as

$$z_k = \prod_{i=1}^j (2 - 2\cos(x - x_i))^{l_i}, \quad k = \sum_{i=1}^j l_i.$$

Hence  $f/z_k$  and  $z_k/f$  have to be bounded and the thesis of part 1 is proved.

*Remark 2.* The preconditioner  $A_n(z_k)$  was first proposed by R. Chan in [7], in which the proof of the related statement of part 1 was also given under the assumption of strong regularity of  $f$  ( $f$  being  $2q$  times continuously differentiable with  $q = \max l_i$ ). The proof under the weaker hypothesis of continuity can be found in [16], while in [30] it is possible to derive the statement under the full

general hypothesis that  $f$  belongs to  $L^1$ ,  $f$  has only “essential” zeros of even order and is essentially nonnegative. For the definition of “essential” zero of a Lebesgue integrable function see [30].

Now we prove statement 2; from the hypothesis  $r < f/g < R$  we deduce that

$$z < \frac{g}{f} < Z, \quad z = \frac{1}{R}, \quad Z = \frac{1}{r}.$$

Consequently, for any positive  $\alpha$  we have

$$(4) \quad \alpha z - 1 < \frac{\alpha g}{f} - 1 = \frac{\alpha g - f}{f} < \alpha Z - 1.$$

Now we choose  $\alpha$  such that  $0 < \alpha Z - 1 = -(\alpha z - 1)$ , i.e., solving the linear equation, we find  $\alpha = 2/(Z + z)$ . Therefore, setting

$$h_{\alpha g} = \left\| \frac{f - \alpha g}{f} \right\|_{\infty}$$

we have from (4)

$$h_{\alpha g} = \alpha Z - 1 = \frac{Z - z}{Z + z} < 1.$$

For the third part it is sufficient to observe that, when  $l \geq k$  we can construct  $z_k \in \mathbf{P}_l$  such that  $r_k < f/z_k < R_k$  and consequently from the former part of the theorem we find  $h_{\alpha z_k} < 1$  and

$$h^* = \min_{g \in \mathbf{P}_l} \left\| \frac{f - g}{f} \right\|_{\infty} \leq \left\| \frac{f - \alpha z_k}{f} \right\|_{\infty} < 1. \quad \square$$

Now define  $g^*$  the polynomial in [11] such that  $h^* = \left\| \frac{f - g^*}{f} \right\|_{\infty}$ ,  $g^A, g^B$  the polynomials shown in the preceding section and scaled suitably in light of the second part of Theorem 4.1, and  $h^A$  and  $h^B$ , respectively, the relative Chebyshev errors.

If  $l \geq k$ , then clearly  $h^* \leq h^A \leq h^B < 1$  so, in view of Theorem 2.2 we have that the condition numbers of the matrices  $A_n^{-1/2}(g^A)A_n(f)A_n^{-1/2}(g^A)$  and  $A_n^{-1/2}(g^B)A_n(f)A_n^{-1/2}(g^B)$  give two upper bounds for

$$\mu_2 \left( A_n^{-1/2}(g^*)A_n(f)A_n^{-1/2}(g^*) \right) < \frac{1 + h^*}{1 - h^*}.$$

We are ready to estimate the condition numbers of the system  $A_n(f)\mathbf{x} = \mathbf{b}$  preconditioned by means of  $A_n(g^A)$  and  $A_n(g^B)$  respectively. In view of Theorem 2.2 we have that the eigenvalues of  $A_n^{-1}(g^A)A_n(f)$  belong to the open interval  $(r_A, R_A) = (\inf \hat{f}/g_{l-k}, \sup \hat{f}/g_{l-k})$ , where  $\hat{f} = f/z_k$  is continuous and positive and  $g_{l-k}$  is the best Chebyshev approximation of  $\hat{f}$ , i.e., we have

$$E_{l-k}(\hat{f}) = \min_{g \in \mathbf{P}_{l-k}} \|\hat{f} - g\|_{\infty} = \|\hat{f} - g_{l-k}\|_{\infty}.$$

Hence

$$\begin{aligned} r_A &\geq \frac{\hat{f}}{\hat{f} + E_{l-k}(\hat{f})} = 1 - \frac{E_{l-k}(\hat{f})}{\hat{f} + E_{l-k}(\hat{f})} \geq 1 - \frac{E_{l-k}(\hat{f})}{\hat{m} + E_{l-k}(\hat{f})}, \\ R_A &\leq \frac{\hat{f}}{\hat{f} - E_{l-k}(\hat{f})} = 1 + \frac{E_{l-k}(\hat{f})}{\hat{f} - E_{l-k}(\hat{f})} \leq 1 + \frac{E_{l-k}(\hat{f})}{\hat{m} - E_{l-k}(\hat{f})}, \end{aligned}$$

$\hat{m}$  being the minimum of  $\hat{f}$  in  $[-\pi, \pi]$ ,  $\hat{m} - E_{l-k}(\hat{f})$  being a positive constant and

$$\frac{E_{l-k}(\hat{f})}{\hat{m} + E_{l-k}(\hat{f})}$$

being less than 1, for  $l - k$  large enough.

Therefore

$$\mu^A = \mu_2 \left( A_n^{-1/2}(g^A) A_n(f) A_n^{-1/2}(g^A) \right) < \frac{1 + \frac{E_{l-k}(\hat{f})}{\hat{m} - E_{l-k}(\hat{f})}}{1 - \frac{E_{l-k}(\hat{f})}{\hat{m} + E_{l-k}(\hat{f})}}.$$

For the second preconditioner  $A_n(g^B)$ ,  $g^B = z_k g_{l-k}$  we recall that  $g_{l-k}$  is the first kind Chebyshev interpolant of  $\hat{f}$  and so, by using the standard error estimate of Powell [25], it follows that

$$\|\hat{f} - g_{l-k}\|_\infty \leq c E_{l-k}(\hat{f}) \log(l - k), \quad c \sim \frac{2}{\pi}.$$

Recalling that  $f/g^B = \hat{f}/g_{l-k}$ , in view of Theorem 2.2 we state that the eigenvalues  $A_n^{-1}(g^B) A_n(f)$  lie in the open interval

$$(r_B, R_B) = (\inf \hat{f}/g_{l-k}, \sup \hat{f}/g_{l-k}).$$

Therefore using the same argument used in analyzing the former preconditioner we obtain

$$\mu^B = \mu_2 \left( A_n^{-1/2}(g^B) A_n(f) A_n^{-1/2}(g^B) \right) < \frac{1 + \frac{F_{l-k}(\hat{f})}{\hat{m} - F_{l-k}(\hat{f})}}{1 - \frac{F_{l-k}(\hat{f})}{\hat{m} + F_{l-k}(\hat{f})}},$$

where  $F_{l-k}(\hat{f}) = c E_{l-k}(\hat{f}) \log(l - k)$ .

Now, to conclude we want to recall a classical result of approximation which allows us to estimate more precisely  $\mu^A$  and  $\mu^B$ .

**Theorem 4.2** ([22]). *If  $f \in C^p[-\pi, \pi]$  and if  $\omega(f^{(p)}; \delta)$  indicates the modulus of continuity of  $f^{(p)}$ , then*

$$E_m(f) \leq d_p \omega \left( f^{(p)}; \frac{1}{m} \right) \frac{1}{m^p},$$

where  $d_p$  is a well-known constant. (Jackson proved this result with a somewhat large constant  $d_p = c^p$ ,  $c < 100$ ; in [23] we find a better estimate of  $c$ , that is,  $c = 1 + \pi^2/2$ .) In particular, if  $f^{(p)} \in \text{Lip}_M^\alpha$ ,  $\alpha \in [0, 1]$ ,  $M > 0$ , (i.e.,  $\forall x_1, x_2 \in [-\pi, \pi]$ ,  $|f^{(p)}(x_1) - f^{(p)}(x_2)| \leq M|x_1 - x_2|^\alpha$ ) we have

$$E_m(f) \leq d_p \frac{1}{m^{p+\alpha}}.$$

Hence

$$\mu^A < U^A(l) = \frac{1 + \frac{d_{l-k} \omega(\hat{f}^{(p)}; \frac{1}{l-k}) / (l-k)^p}{\hat{m} - d_{l-k} \omega(\hat{f}^{(p)}; \frac{1}{l-k}) / (l-k)^p}}{1 - \frac{d_{l-k} \omega(\hat{f}^{(p)}; \frac{1}{l-k}) / (l-k)^p}{\hat{m} + d_{l-k} \omega(\hat{f}^{(p)}; \frac{1}{l-k}) / (l-k)^p}},$$

where  $U^A(l)$  can be a sharp estimate of  $\mu_2^A$  for a large class of functions for which the estimate of Jackson (see the last theorem) is sharp.



5. SOME REMARKS ABOUT THE COMPUTATIONAL COSTS

In this section first we discuss the computational cost required by each iteration of the considered PCG methods and then we discuss the cost of the determination of the trigonometric polynomials  $g^*$ ,  $g^A$  and  $g^B$ . For the first step we consider two possible strategies: the first one based on classical band solvers [19] and the second one based on a particular algebraic multigrid method devised for symmetric positive definite Toeplitz matrices [17]. By following the first idea, we have that the solution of a system  $A_n(g^A)\mathbf{y} = \mathbf{c}$  needs about  $q_1 l^2 n$  arithmetic ops while the second method requires  $q_2 l n$  ops, that is, the cost is linear both with respect to the dimension and to the bandwidth. Since the multiplication of  $A_n(f)$  by a vector uses about  $q_3 n \log n$  arithmetic ops ( $q_i$  suitable and known constants) we can state that the number of iterations to reach the solution within a preassigned accuracy  $\epsilon$ , is proportional to

$$\text{Cost}(l, n, \epsilon) = nU^A(l)(X + q_2 \log n) \log\left(\frac{1}{\epsilon}\right)$$

where  $X$  is equal to  $q_1 l^2$  or to  $q_2 l$ .

Therefore, observing that  $U^A(l)$ ,  $U^*(l)$  and  $U^B(l)$  are nonincreasing functions of  $l$  independent of  $n$ , we obtain that a total cost of  $O(n \log n)$  ops is reached for any  $l \in [k, l_{\max}]$  where  $l_{\max} = O(\log^{1/2} n)$  if we use classical band solvers or  $l_{\max} = O(\log n)$  if we use the quoted algebraic multigrid method. Finally, the global optimal value  $l_{opt}$  can be estimated by minimizing analytically or numerically the function  $\text{Cost}(l, n, \epsilon)$  with respect to the variable  $l$ . For instance, by considering Table 4 in section 7, if we use the multigrid strategy [17] for the solution of the banded preconditioning systems, then the optimal value of the halfbandwidth is  $l = 6$ .

**5.1. The computation of the “generating” functions  $g^*$ ,  $g^A$  and  $g^B$ .** We preliminarily observe that the calculation of  $g^*$ , which involves a modified version of the Remez algorithm [33, 12], and the calculation of  $g^A$ , which requires the standard Remez algorithm [26], are from a computational point of view substantially equivalent. In both of the cases we use an iterative method which has a complex structure (see [12]) and for which we cannot give a theoretical bound of the number of iterations required for the computation within a preassigned accuracy (see [26, 11, 12, 33]. However, in practice in [33, 11], the authors observe a total arithmetic amount of about  $O(l^3)$  ops where  $l$  is the degree of the trigonometric polynomials  $g^*$  or, equivalently,  $g^B$ .

On the other hand, the calculation of  $g^*$  is delicate. By referring to the paper [12], in step 1, the modified Remez algorithm has to solve a linear system in which the elements of the  $j$ -th column ( $j = 2, \dots, l + 1$ ) are given by  $\phi_{j-1}(x_s^{(i)})$  where  $\phi_k(x) = \frac{\cos(k-1)x}{f(x)}$  and  $x_s^{(i)}$ ,  $s = 1, \dots, l + 1$ , are points calculated at the  $i$ -th step of the algorithm. If some of the values approach a zero of  $f$  and, in particular, when  $f$  has zeros of high order, it may happen that the considered linear system becomes very ill-conditioned.

The calculation of the coefficients of the polynomial  $g^B$  is actually very simple and direct.

Let us consider a function  $h$  defined on  $[-1, 1]$  and a positive integer  $m$ , then we consider the trigonometric polynomial of degree  $m$  interpolating  $h$  at zeros of

the  $m + 1$ -th Chebyshev polynomial of the first kind. A classical trigonometric representation is given by

$$b_0 + 2 \sum_{j=1}^m b_j \cos(jx)$$

where the coefficients  $b_i$  are easily expressed as follows:

$$b_j = \frac{1}{m+1} \sum_{k=1}^m h(\cos(\theta_k)) \cos(j\theta_k), \quad \theta_k = \frac{(2k+1)\pi}{2(m+1)}.$$

Therefore, by using about  $2m^2$  arithmetic ops, we may obtain all the coefficients  $b_j$ . On the other hand, the previous expression can be suitably manipulated in the following way:

$$b_j = \frac{1}{2(m+1)} \sum_{k=1}^m h(\cos(\theta_k)) (e^{ij\theta_k} + e^{-ij\theta_k}), \quad \theta_j = \frac{j\pi}{m+1} + \frac{\pi}{2(m+1)}.$$

By calling  $\phi_j = \frac{j\pi}{m+1}$  we have that  $b_j = z_j + \bar{z}_j$  where

$$z_j = \frac{1}{2(m+1)} e^{ij \frac{\pi}{2(m+1)}} \left[ \sum_{k=1}^m h(\cos(\theta_k)) e^{ij\phi_k} \right]$$

and the expression in the squared brackets can be calculated by using one FFT of order  $2(m+1)$ . Therefore the calculation of all the  $b_j$  can be done in a stable way within a cost of  $cm \log m$  where  $c$  is a suitable, small and known constant value [34].

Obviously, by using the approximation theory tools we can define other good preconditioners whose generating functions can be efficiently calculated by means of FFTs algorithms. For instance, if we consider the least squares polynomial  $g$  of degree  $m$  [25] and we approximate its coefficients by using the repeated trapezoidal rule with  $m+2$  uniformly distributed knots, we obtain some “cosine” summations which can be calculated by using fast cosine transforms [1]. We recall that the order of approximation of the Chebyshev least squares polynomial is the same as that of the polynomial interpolating at the Chebyshev zeros [25].

Finally, if we are interested in parallel computation (for instance the PRAM model), it is useful to recall that FFT algorithms are well parallelizable procedures and perform  $O(\log m)$  parallel steps if  $m$  is the order of the Fourier transform.

Since the band solvers [19] are inherently sequential we cannot use them for the solution of the preconditioning systems. However, in the recent literature we may find alternative techniques: for an efficient parallel solution of generic band systems see [35], while a good survey about fast parallel methods for band-Toeplitz systems can be found in [4]. Since the trigonometric polynomials  $g$  are nonnegative, it follows that the matrices  $A_n(g)$  are also positive definite; therefore we can alternatively apply the multigrid technique developed in [17].

## 6. SUPERLINEAR PCG METHODS

In this section we discuss a strategy in order to devise superlinear PCG methods having the optimal cost of  $O(n \log n)$  ops. We say that this cost is *optimal* because  $O(n \log n)$  ops is, asymptotically, the cost of a FFT and, therefore, the cost of a product between a (dense) Toeplitz matrix and a generic vector. Since in the

implementation of a PCG method we have to calculate, at each iteration, a few of these products, it seems evident that the asymptotical cost of  $O(n \log n)$  ops is minimal with respect to this class of linear algebra problems.

Then, if we use a classic band solver [19], then the maximal halfbandwidth we can admit is  $l_{\max} = O(\log^{1/2} n)$ , while, if we choose an algebraic multigrid method as in [17] we can use a maximal halfbandwidth equal to  $l_{\max} = O(\log n)$ . So, by calling  $t_n = l_{\max}^{-1}$ , we find that  $U^A(l_{\max})$  is asymptotic to

$$1 + O\left(t_n^p \omega\left(\hat{f}^{(p)}; t_n^p\right)\right),$$

$U^*(l_{\max})$  is less than  $U^A(l_{\max})$  and  $U^B(l_{\max})$  is asymptotic to

$$1 + O\left(\log(l_{\max}) t_n^p \omega\left(\hat{f}^{(p)}; t_n^p\right)\right).$$

In this way we have found PCG methods with a cost of  $O(n \log n)$  ops and having a superlinear rate of convergence due to the cluster around 1 observed for the spectrum of the preconditioned matrices  $A_n^{-1}(g^A)A_n(f)$  and  $A_n^{-1}(g^*)A_n(f)$ . In fact, the functions  $U^*(l_{\max})$  and  $U^A(l_{\max})$  tend to 1 as the dimension  $n$  tends to infinity.

Of course, analogous considerations can be done naturally for the PCG method associated with the preconditioner  $A_n(g^B)$ . In this last case the “clustering” property is weakly deteriorated according to the quantity

$$\log(l_{\max}(n)) = O(\log(\log n)).$$

This deterioration becomes considerable only when the function  $\hat{f}$  is very “irregular”. Actually, it is sufficient that the modulus of continuity  $\omega(f; \delta)$  of  $f$  is a “small  $o$ ” of  $\frac{1}{|\log \delta|}$  in order to have

$$\lim_{n \rightarrow \infty} U^B(l_{\max}(n)) = 1.$$

For instance, the very weak assumption that  $\hat{f}$  belongs to the class  $\text{Lip}_M^\alpha$  for some positive value  $\alpha$  is enough in order to obtain the preceding relation.

If we want to use in practice these superlinear PCG methods, we have to point out that the considered trigonometric polynomials  $g^*$ ,  $g^A$  and  $g^B$  have a degree which grows logarithmically as a function of  $n$ . This means that, for any dimension  $n$ , we must calculate the generating function of the preconditioner. Unfortunately, for the first two strategies we cannot exhibit a theoretical bound for the related arithmetic cost; for the third strategy, on the other hand, we may say that the cost is of  $O(l_{\max}(n) \log l_{\max}(n))$  (see the previous section). Since  $l_{\max}(n)$  is bounded by  $c_1 \log^{1/2}(n)$  or  $c_2 \log(n)$  we have that the total cost for determining the coefficients of  $g^B$  is well dominated by the asymptotical (and also practical) cost of a generic iteration of the associated PCG method. Therefore, as shown in the last section, we can really use this superlinear technique.

Observe that good clustering properties for the preconditioned matrix have been proved also in [8], where by using band-Toeplitz + circulant preconditioners the authors may handle the non-Hermitian case. By comparing the two approaches, we notice that the techniques proposed in this paper do not seem to be interesting for non-Hermitian problems, because, for instance, the theoretical results as in Theorem 2.2 and Theorem 2.3 are strictly related to the symmetric case. However, we stress that in our case we can easily deduce a superlinear convergence by using

the previous limit relations and the tools developed by Axelsson and Lindskög [3], while in [8] no rigorous proof of superlinear convergence is given.

## 7. NUMERICAL RESULTS

In this section, we compare the convergence rate of the band-Toeplitz preconditioner (strategy **A**), with the optimal band-Toeplitz preconditioner [11] and with the optimal circulant preconditioner [13] on three different generating functions having zeros. They are  $(x - 1)^2(x + 1)^2$ ,  $1 - e^{-x^2}$  and  $x^4$  and are associated to ill-conditioned matrices  $A_n$  having Euclidean condition numbers equal to  $O(n^2)$ ,  $O(n^2)$  and  $O(n^4)$  respectively (see for instance [30, 31]). The matrices  $A_n$  are formed by evaluating the Fourier coefficients of the generating functions by using FFTs (see [11]). In the tests considered, the component of the vector  $\mathbf{b}$  on the right-hand side of the system  $A_n \mathbf{x} = \mathbf{b}$  are all equal to one, the zero vector is the initial guess and the stopping criterion is  $\|r_q\|_2 / \|r_0\|_2 \leq 10^{-7}$ , where  $r_q$  is the residual vector after  $q$  iterations. All computations were performed using Matlab.

In the subsequent tables,  $I$  denotes that no preconditioning is used,  $C$  is the T. Chan optimal circulant preconditioner [13],  $B_{n,l}^*$  is the optimal band-Toeplitz preconditioner [11] and  $B_{n,l}^B$  is the band-Toeplitz preconditioner defined according to the strategy **B**; here  $l$  denotes the halfbandwidth of the band preconditioners.

We do not make explicit comparison with the preconditioner related to the strategy **A** because the associated PCG method has, by virtue of the relation  $\mu^* < \mu^A < \mu^B$ , a convergence speed between the R. Chan, P. Tang one and the “**B**” one.

We observe that the “optimal” and the “**B**” band-Toeplitz PCG methods perform, substantially, in the same way, but the second one is much more economical with respect to the computation of the related generating function. This fact is not so considerable when the bandwidth is fixed, but it becomes crucial in order to increase  $l$ , say, as  $\log n$ . Actually, in this case, for any dimension  $n$ , it is not expensive to calculate a different preconditioner  $A_n(g^A(l))$ , since the related cost  $O(\log n \log(\log n))$  is strongly dominated by the cost  $O(n \log n)$  of each PCG iteration.

Finally, the reduction of the number of required iterations, as the dimension increases, shown in Table 4 gives numerical evidence of the superlinear convergence claimed in section 6. We stress that the exceptional convergence behaviour of the PCG algorithm related to  $B_{n,6}^B$  is explained by the good approximation properties of the first-kind Chebyshev interpolation: to have a practical measure of this,

TABLE 1.  $f(x) = (x^2 - 1)^2$

$n$	$I$	$C$	$B_{n,3}^* = B_{n,3}^B$	$B_{n,4}^* \ B_{n,4}^B$	$B_{n,5}^* \ B_{n,5}^B$	$B_{n,6}^* \ B_{n,6}^B$
16	11	9	9	9 7	8 6	7 6
32	27	14	13	11 9	9 7	7 6
64	74	17	16	11 10	8 8	7 7
128	193	22	18	11 11	8 8	7 7
256	465	28	19	11 11	8 9	7 7
512	> 1000	34	19	11 11	8 8	7 7

it is sufficient to notice that the reduction of the condition number from  $A_n$  to  $(B_{n,6}^B)^{-1}A_n$ , for  $n = 512$ , is from  $2.7 * 10^4$  to  $1 + 5 * 10^{-4}$ .

TABLE 2.  $f(x) = 1 - e^{-x^2}$

$n$	$I$	$C$	$B_{n,2}^* = B_{n,2}^B$	$B_{n,3}^* \ B_{n,3}^B$	$B_{n,4}^* \ B_{n,4}^B$	$B_{n,5}^* \ B_{n,5}^B$
16	9	6	9	7 8	4 4	3 3
32	14	7	15	7 8	5 5	3 3
64	24	8	17	8 9	5 5	3 3
128	42	10	17	8 9	5 5	3 3
256	77	13	17	8 9	5 5	3 3
512	143	17	17	8 9	5 5	3 3

TABLE 3.  $f(x) = x^4$

$n$	$I$	$C$	$B_{n,3}^* = B_{n,3}^B$	$B_{n,4}^* \ B_{n,4}^B$	$B_{n,5}^* \ B_{n,5}^B$	$B_{n,6}^* \ B_{n,6}^B$
16	12	10	9	9 8	9 7	7 6
32	34	16	15	10 10	11 8	9 7
64	119	26	21	13 12	11 10	9 8
128	587	77	24	15 15	12 11	10 10
256	> 1000	179	27	16 16	12 13	10 10
512	> 1000	406	29	16 16	13 13	10 11

TABLE 4.  $f(x) = 1 - e^{-x^2}$ , superlinear PCG Prec =  $B_{n,l(n)}^B$ ,  $l(n) = \log(n) - 2$

$n$	16	32	64	128	256	512
$l(n)$	2	3	4	5	6	7
Iter	9	7	5	3	2	2

ACKNOWLEDGEMENT

It is a pleasure for me to thank Professor D. Bini and Professor M. Capovani for their suggestions and for inspiring me to do mathematical research.

REFERENCES

1. N. Ahmed, T. Natarajan, K. Rao, "Discrete cosine transforms", *IEEE Trans. Comp.*, **23** (1974), 90-93. MR **50**:9025
2. O. Axelsson, V. Barker, *Finite Element Solution of Boundary Value Problems, Theory and Computation*, Academic press Inc., New York 1984. MR **85m**:65116
3. O. Axelsson, G. Lindskog, "The rate of convergence of the preconditioned conjugate gradient method", *Num. Math.*, **52** (1986), 499-523. MR **88a**:65037b
4. D. Bini, "Matrix structure in parallel matrix computation", *Calcolo*, **25** (1988), pp. 37-51. MR **91g**:65322
5. D. Bini, M. Capovani, "Spectral and computational properties of band symmetric Toeplitz matrices", *Linear Algebra Appl.*, **52/53** (1983), 99-126. MR **85k**:15008

6. D. Bini, P. Favati, "On a matrix algebra related to the discrete Hartley transform", *SIAM J. Matrix Anal. Appl.*, **14** (1993), 500–507. MR **94h**:65026
7. R.H. Chan, "Toeplitz preconditioners for Toeplitz systems with nonnegative generating functions", *IMA J. Numer. Anal.*, **11** (1991), 333–345. MR **92f**:65041
8. R.H. Chan, W. Ching, "Toeplitz–circulant preconditioners for Toeplitz systems and their applications to queueing network with batch arrivals", *SIAM J. Sci. Comp.*, **17** (1996), 762–772. CMP 96:11
9. R.H. Chan, Q. Chang, H. Sun, "Multigrid method for ill-conditioned symmetric Toeplitz systems", *personal communication*.
10. R.H. Chan, G. Strang, "Toeplitz equations by conjugate gradients with circulant preconditioner", *SIAM J. Sci. Stat. Comp.*, **10** (1989), 104–119. MR **90d**:65069
11. R.H. Chan, P. Tang, "Fast band–Toeplitz preconditioners for Hermitian Toeplitz systems", *SIAM J. Sci. Comp.*, **15** (1994), 164–171. MR **94j**:65043
12. R.H. Chan, P. Tang, "Constrained minimax approximation and optimal preconditioners for Toeplitz matrices", *Numer. Alg.*, **5** (1993), 353–364. CMP 94:07
13. T.F. Chan, "An optimal circulant preconditioner for Toeplitz systems", *SIAM J. Sci. Stat. Comp.*, **9** (1988), 766–771. MR **89e**:65046
14. P. Davis, *Circulant Matrices*. John Wiley and Sons, New York 1979. MR **81a**:15003
15. F. Di Benedetto, "Analysis of preconditioning techniques for ill-conditioned Toeplitz matrices", *SIAM J. Sci. Comp.*, **16** (1995), 682–697. MR **95m**:65082
16. F. Di Benedetto, G. Fiorentino, S. Serra, "C.G. Preconditioning for Toeplitz Matrices", *Comp. Math. Applic.*, **25** (1993), 35–45. MR **93h**:65063
17. G. Fiorentino, S. Serra, "Multigrid methods for Toeplitz matrices", *Calcolo*, **28** (1991), pp. 283–305. MR **94c**:65039
18. I. Gohberg, I. Feldman, *Convolution Equations and Projection Methods for Their Solution*, Transaction of Mathematical Monographs, **41**, American Mathematical Society, Providence, Rhode Island 1974. MR **50**:8149
19. G.H. Golub, C.F. Van Loan, *Matrix Computations*. The Johns Hopkins University Press, Baltimore 1983. MR **85h**:65063
20. U. Grenander, G. Szegő, *Toeplitz Forms and Their Applications*. Second Edition, Chelsea, New York, 1984. MR **88b**:42031
21. I.S. Iokhvidov, *Hankel and Toeplitz Forms: Algebraic Theory*. Birkhäuser, Boston, 1982. MR **83k**:15021
22. D. Jackson, *The Theory of Approximation*. American Mathematical Society, New York, 1930.
23. G. Meinardus, *Approximation of Functions: Theory and Numerical Methods*. Springer–Verlag, Berlin, 1967. MR **36**:571
24. A. Oppenheim, *Applications of Digital Signal Processing*. Prentice–Hall, Englewood Cliffs, 1978.
25. M.J.D. Powell, "On the maximum errors of polynomial approximation defined by interpolation and least squares criteria", *Comput. J.*, **9** (1966), 404–407. MR **34**:8616
26. E.J. Remes, "Sur le calcul effectif des polynomes d'approximation de Tchebichef" *Compt. Rend. Acad. Sci. Paris*, **199** (1934), 337–340.
27. S. Serra, "Preconditioning strategies for asymptotically ill-conditioned block Toeplitz systems", *BIT*, **34** (1994), pp. 579–594.
28. S. Serra, "Conditioning and solution of Hermitian (block) Toeplitz systems by means of preconditioned conjugate gradient methods", *Proc. in Advanced Signal Processing Algorithms, Architectures, and Implementations - SPIE conference*, F. Luk Ed., San Diego (CA), July 1995, pp. 326–337.
29. S. Serra, "Preconditioning strategies for Hermitian Toeplitz systems with nondefinite generating functions", *SIAM J. Matr. Anal. Appl.*, **17** (1996), pp. 1007–1019.
30. S. Serra, "On the extreme eigenvalues of Hermitian (block) Toeplitz matrices", *Linear Algebra Appl.*, to appear.
31. S. Serra, "On the extreme spectral properties of Toeplitz matrices generated by  $L^1$  functions with several minima (maxima)", *BIT*, **36** (1996), 135–142.
32. S. Serra, "Asymptotic results on the spectra of preconditioned Toeplitz matrices and some applications", *TR nr. 9 University of Calabria*, (1995).
33. P. Tang, "A fast algorithm for linear complex Chebyshev approximation", *Math. Comp.*, **51** (1988), 721–739. MR **89j**:30054

34. C. Van Loan, *Computational Frameworks for the Fast Fourier Transform*. SIAM, Philadelphia, 1992. MR **93a**:65186
35. S. Wright, "Parallel algorithms for banded linear systems", *SIAM J. Sci. Stat. Comp.*, **12** (1991), 824–842. MR **92a**:65096

DIPARTIMENTO DI INFORMATICA, UNIVERSITÀ DI PISA, CORSO ITALIA 40, 56100 PISA (ITALY)  
*E-mail address:* `serra@morse.dm.unipi.it`