

## CIRCULANT PRECONDITIONERS FOR TOEPLITZ MATRICES WITH POSITIVE CONTINUOUS GENERATING FUNCTIONS

RAYMOND H. CHAN AND MAN-CHUNG YEUNG

**ABSTRACT.** We consider the solution of  $n$ -by- $n$  Toeplitz systems  $A_n x = b$  by the preconditioned conjugate gradient method. The preconditioner  $C_n$  is the circulant matrix that minimizes  $\|B_n - A_n\|_F$  over all circulant matrices  $B_n$ . We show that if the generating function  $f$  is a positive  $2\pi$ -periodic continuous function, then the spectrum of the preconditioned system  $C_n^{-1} A_n$  will be clustered around one. In particular, if the preconditioned conjugate gradient method is applied to solve the preconditioned system, the convergence rate is superlinear.

### 1. INTRODUCTION

In this paper, we discuss the solution of Toeplitz systems  $A_n x = b$  by the preconditioned conjugate gradient method. The idea of using the method with circulant preconditioners for solving symmetric positive definite Toeplitz systems was first proposed by Strang [11]. The number of operations per iteration is  $O(n \log n)$  as circulant systems can be solved efficiently by the Fast Fourier Transform and the matrix-vector multiplication  $A_n x$  can also be computed by Fast Fourier Transform by first embedding  $A_n$  into a  $2n$ -by- $2n$  circulant matrix.

Several other circulant preconditioners have been proposed and analyzed since then; see for instance, R. Chan and Strang [1], R. Chan [2, 3], R. Chan, Jin, and Yeung [5], T. Chan [6], T. Ku and Kuo [9], and E. Tyrtshnikov [12]. It has been shown in these papers that if the diagonals of the Toeplitz matrix  $A_n$  are Fourier coefficients of a positive function in the Wiener class, then the spectrum of the preconditioned system will be clustered around one. It follows that the preconditioned conjugate gradient method, when applied to solve the preconditioned system, converges superlinearly. Hence the number of iterations required for convergence is independent of the size of the matrix  $A_n$ . In particular, the system  $A_n x = b$  can be solved in  $O(n \log n)$  operations.

It has also been proved in R. Chan [3] and R. Chan, Jin, and Yeung [5] that, under the same Wiener class assumption, the preconditioned systems with preconditioners proposed by either R. Chan [3], T. Chan [6], Strang [11], or

---

Received October 25, 1990.

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

*Key words and phrases.* Toeplitz matrix, circulant matrix, preconditioned conjugate gradient method, generating function.

Tyrtyshnikov [12] have spectra that are asymptotically the same. In particular, all these preconditioned systems converge at the same rate for large  $n$ . However, Tyrtyshnikov [12] showed that for general positive definite Toeplitz matrix  $A_n$ , the corresponding T. Chan's and Tyrtyshnikov's preconditioners are also positive definite, but R. Chan's and Strang's preconditioners are in general not. We note that it requires  $6n \log n + O(n)$  operations to generate Tyrtyshnikov's preconditioner (see R. Chan, Jin, and Yeung [4]), but only  $3n/2$  operations to form T. Chan's preconditioner. Thus T. Chan's preconditioner is the best choice for Toeplitz systems that satisfy the Wiener class assumption. In the following, we therefore only consider T. Chan's preconditioner.

The main result in this paper is to extend the above-mentioned superlinear convergence result from the Wiener class of functions to the class of  $2\pi$ -periodic continuous functions. More precisely, we will show that if the diagonals of  $A_n$  are Fourier coefficients of a positive  $2\pi$ -periodic continuous function, and if  $C_n$  is the corresponding T. Chan's preconditioner for  $A_n$ , then the spectrum of  $C_n^{-1}A_n$  will be clustered around one and the preconditioned system converges superlinearly. The outline of the paper is as follows. In §2, we discuss some of the properties of the spectra of  $A_n$  and  $C_n$ . In §3, Jackson's formula from approximation theory is introduced and used to derive our main theorem. Finally, numerical examples are given in §4.

## 2. THE SPECTRA OF $C_n$ AND $A_n$

For simplicity, we denote by  $\mathcal{E}_{2\pi}$  the Banach space of all  $2\pi$ -periodic continuous real-valued functions equipped with the supremum norm  $\|\cdot\|_\infty$ . For all  $f \in \mathcal{E}_{2\pi}$ , let

$$a_k(f) = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(\theta) e^{-ik\theta} d\theta, \quad k = 0, \pm 1, \pm 2, \dots,$$

be the Fourier coefficients of  $f$ . Since  $f$  is real-valued,  $a_{-k} = \bar{a}_k$  for all integers  $k$ . Let  $A_n(f)$  be the  $n$ -by- $n$  Hermitian Toeplitz matrix with the  $(j, k)$ th entry given by  $a_{j-k}(f)$ . The function  $f$  is called the generating function of the matrices  $A_n(f)$ . The following lemma gives the relation between  $f$  and the spectrum  $\sigma(A_n(f))$  of  $A_n(f)$ . The proof is given in Grenander and Szegő [8, p. 65].

**Lemma 1.** *Let  $f \in \mathcal{E}_{2\pi}$ , with the minimum and maximum values given by  $f_{\min}$  and  $f_{\max}$ , respectively. Then  $\sigma(A_n(f)) \subseteq [f_{\min}, f_{\max}]$ . In particular, we have*

$$(1) \quad \|A_n(f)\|_2 \leq \|f\|_\infty.$$

Let  $C_n(f)$  be the  $n$ -by- $n$  circulant preconditioner of  $A_n(f)$  as defined in T. Chan [6], i.e.,  $C_n(f)$  is the minimizer of  $\|B_n - A_n(f)\|_F$  over all  $n$ -by- $n$  circulant matrices  $B_n$ . We note that the diagonals  $c_k$  of  $C_n$  can be obtained by averaging the corresponding diagonals of  $A_n$ , with the diagonals of  $A_n$  being extended to length  $n$  by a wrap-around. More precisely, the  $c_k$  are given by

$$(2) \quad c_k = \begin{cases} ((n-k)a_k + ka_{k-n})/n, & 0 \leq k < n, \\ \bar{c}_{-k}, & 0 < -k < n. \end{cases}$$

The following lemma gives the relationship between the spectra of  $C_n(f)$  and  $A_n(f)$ , and its proof can be found in R. Chan, Jin, and Yeung [4]; see also Tyrtyshnikov [12].

**Lemma 2.** Let  $A_n$  be an arbitrary  $n$ -by- $n$  Hermitian matrix and  $C_n$  be the optimal preconditioner given in (2). Then  $C_n$  is Hermitian and

$$\lambda_{\min}(A_n) \leq \lambda_{\min}(C_n) \leq \lambda_{\max}(C_n) \leq \lambda_{\max}(A_n),$$

where  $\lambda_{\max}(\cdot)$  and  $\lambda_{\min}(\cdot)$  denote the largest and the smallest eigenvalues respectively. In particular, if  $A_n$  is positive definite, then  $C_n$  is also positive definite.

Combining Lemmas 1 and 2, we have the following immediate corollary.

**Corollary 1.** Let  $f \in \mathcal{E}_{2\pi}$ ; then  $\sigma(C_n(f)) \subseteq [f_{\min}, f_{\max}]$ . In particular, we have

$$(3) \quad \|C_n(f)\|_2 \leq \|f\|_{\infty}.$$

Moreover, if  $f$  is positive, then for all  $n > 0$ ,  $A_n(f)$  and  $C_n(f)$  are positive definite and  $C_n(f)$  and  $C_n^{-1}(f)$  are uniformly bounded in the  $l_2$ -norm.

The proof of the next lemma is given in R. Chan [3]. It basically states that the spectrum of  $A_n(f) - C_n(f)$  is clustered around zero if  $f$  is in the Wiener class. We remark that a function  $f$  is in the Wiener class if its Fourier coefficients are absolutely summable, i.e.,

$$\sum_{k=-\infty}^{\infty} |a_k(f)| < \infty.$$

**Lemma 3.** Let  $f$  be a function in the Wiener class; then for all  $\varepsilon > 0$ , there exist  $N, M > 0$  such that for all  $n > N$ , at most  $M$  eigenvalues of  $A_n(f) - C_n(f)$  have absolute value larger than  $\varepsilon$ .

The main result of this paper is to extend the result in Lemma 3 from the Wiener class of functions to  $\mathcal{E}_{2\pi}$ . We first note that if  $f$  is in the Wiener class, then  $f \in \mathcal{E}_{2\pi}$ . In fact, if the Fourier coefficients of  $f$  are absolutely summable, then the Fourier series of  $f$  is a well-defined function in  $\mathcal{E}_{2\pi}$ . Moreover, by the Weierstrass test (see Rudin [10, p. 148]), this Fourier series converges uniformly to  $f$  on  $[-\pi, \pi]$ . Hence,  $f$  itself must be in  $\mathcal{E}_{2\pi}$ . On the other hand, the Hardy-Littlewood series given by

$$(4) \quad H(\theta) = \sum_{k=1}^{\infty} \left\{ \frac{e^{ik \log k}}{k} e^{ik\theta} + \frac{e^{-ik \log k}}{k} e^{-ik\theta} \right\}$$

is a classical example of a function which is in  $\mathcal{E}_{2\pi}$  but not in the Wiener class (see Zygmund [14, p. 197]). Thus, the Wiener class of functions is a proper subset of  $\mathcal{E}_{2\pi}$ .

### 3. THE SPECTRUM OF $A_n(f) - C_n(f)$

The idea of our proof is to use the Weierstrass theorem to approximate any function in  $\mathcal{E}_{2\pi}$  by trigonometric polynomials. However, in order to obtain more precise information on the distribution of the eigenvalues, we resort to a stronger form of the Weierstrass theorem, called the Jackson formula, which is given in Lemma 4 below, whose proof can be found in Cheney [7, p. 144]. We first denote the space of all  $n$ th-degree real trigonometric polynomials by  $\mathcal{P}_n$ , i.e.,

$$\mathcal{P}_n = \left\{ \sum_{k=-n}^n b_k e^{ik\theta} \mid b_{-k} = \bar{b}_k, \forall |k| \leq n \right\}.$$

**Lemma 4.** *Let  $f \in \mathcal{E}_{2\pi}$ . Then for all  $n > 0$ ,*

$$\inf_{p_n \in \mathcal{P}_n} \|f - p_n\|_\infty \leq \omega\left(f; \frac{\pi}{n+1}\right),$$

where  $\omega$  is the modulus of continuity of  $f$ , i.e.,

$$\omega(f; \delta) = \sup_{|\theta_1 - \theta_2| \leq \delta} |f(\theta_1) - f(\theta_2)|.$$

Since  $\mathcal{P}_n$  is a finite-dimensional subspace of the normed space  $\mathcal{E}_{2\pi}$ , the infimum of the continuous functional  $\|f - \cdot\|_\infty$  in Lemma 4 above is attained by some polynomials in  $\mathcal{P}_n$ . In particular, we have the following corollary.

**Corollary 2.** *Let  $f \in \mathcal{E}_{2\pi}$ . Then for all  $n > 0$ , there exists a trigonometric polynomial  $p_n \in \mathcal{P}_n$  such that*

$$\|f - p_n\|_\infty \leq \omega\left(f; \frac{\pi}{n+1}\right).$$

The next theorem states that the spectrum of  $A_n(f) - C_n(f)$  is clustered around zero. We prove it by showing that  $A_n(f) - C_n(f)$  can be written as the sum of a low-rank matrix and a small-norm matrix.

**Theorem 1.** *Let  $f \in \mathcal{E}_{2\pi}$ . Then for all  $\varepsilon > 0$ , there exist  $N$  and  $M > 0$  such that for all  $n > N$ , at most  $M$  eigenvalues of  $A_n(f) - C_n(f)$  have absolute values larger than  $\varepsilon$ .*

*Proof.* Let  $f \in \mathcal{E}_{2\pi}$ . Then for any  $\varepsilon > 0$ , there exists a  $\delta > 0$  such that

$$|f(\theta_1) - f(\theta_2)| < \varepsilon \quad \forall |\theta_1 - \theta_2| < \delta.$$

Let  $M = \lceil \pi/\delta \rceil$ ; then

$$\omega\left(f; \frac{\pi}{M+1}\right) \leq \varepsilon.$$

Hence, by Corollary 2, there is a trigonometric polynomial

$$p_M(\theta) = \sum_{k=-M}^M b_k e^{ik\theta}$$

with  $b_{-k} = \bar{b}_k$  such that

$$(5) \quad \|f - p_M\|_\infty \leq \omega\left(f; \frac{\pi}{M+1}\right) \leq \varepsilon.$$

For all  $n > 2M$ , we write

$$(6) \quad \begin{aligned} C_n(f) - A_n(f) &= C_n(f - p_M) - A_n(f - p_M) + C_n(p_M) - A_n(p_M) \\ &= C_n(f - p_M) - A_n(f - p_M) + W_n + U_n, \end{aligned}$$

where by (2) we see that  $W_n$  and  $U_n$  are Hermitian Toeplitz matrices with the first row given by

$$(7) \quad \left(0, -\frac{1}{n}b_{-1}, \dots, -\frac{M}{n}b_{-M}, 0, \dots, 0\right)$$

and

$$(8) \quad \left(0, \dots, 0, \frac{n-M}{n}b_M, \dots, \frac{n-1}{M}b_1\right),$$

respectively. It is clear from (8) that

$$(9) \quad \text{rank } U_n \leq 2M.$$

We will show that the first three terms in the right-hand side of (6) are matrices of small norm. We note that by (1), (3), and (5),

$$(10) \quad \begin{aligned} \|C_n(f - p_M) - A_n(f - p_M)\|_2 &\leq \|C_n(f - p_M)\|_2 + \|A_n(f - p_M)\|_2 \\ &\leq \|f - p_M\|_\infty + \|f - p_M\|_\infty \leq 2\varepsilon. \end{aligned}$$

It remains to estimate  $\|W_n\|_2$ . For all  $|k| \leq M$ , we first note that

$$\begin{aligned} |b_k| &= \left| \frac{1}{2\pi} \int_{-\pi}^{\pi} p_M(t) e^{-ikt} dt \right| \\ &\leq \left| \frac{1}{2\pi} \int_{-\pi}^{\pi} (p_M(t) - f(t)) e^{-ikt} dt \right| + \left| \frac{1}{2\pi} \int_{-\pi}^{\pi} f(t) e^{-ikt} dt \right| \\ &\leq \|f - p_M\|_\infty + f_{\max} \leq \varepsilon + f_{\max}. \end{aligned}$$

Since  $W_n$  is Hermitian, we see from (7) that

$$\begin{aligned} \|W_n\|_2 &\leq \|W_n\|_\infty = 2 \cdot \left( \frac{1}{n} |b_1| + \frac{2}{n} |b_2| + \dots + \frac{M}{n} |b_M| \right) \\ &\leq 2 \cdot \frac{1}{n} \cdot (1 + 2 + \dots + M) (\varepsilon + f_{\max}) \\ &= \frac{1}{n} M(M + 1) (\varepsilon + f_{\max}). \end{aligned}$$

Therefore, if we let

$$\begin{aligned} N &\equiv \max \left\{ M(M + 1) \left( 1 + \frac{f_{\max}}{\varepsilon} \right), 2M \right\} \\ &= M(M + 1) \left( 1 + \frac{f_{\max}}{\varepsilon} \right) = \left[ \frac{\pi}{\delta} \right] \left( \left[ \frac{\pi}{\delta} \right] + 1 \right) \left( 1 + \frac{f_{\max}}{\varepsilon} \right), \end{aligned}$$

then for all  $n \geq N$ , we have  $\|W_n\|_2 \leq \varepsilon$ . Thus, combining this estimate with (9) and (10), we see that for all  $n \geq N$ ,  $C_n(f) - A_n(f)$  is the sum of a matrix of  $l_2$ -norm less than  $3\varepsilon$  and a matrix of rank less than  $2M$ . Hence by using the Cauchy interlace theorem (see, for instance, Wilkinson [13, p. 103]), we see that the matrix  $C_n(f) - A_n(f)$  has at most  $2M$  eigenvalues with absolute values larger than  $3\varepsilon$  whenever  $n > N$ .  $\square$

If  $f$  is Lipschitz continuous, i.e., there exists an  $L$  such that

$$|f(\theta_1) - f(\theta_2)| \leq L|\theta_1 - \theta_2| \quad \forall \theta_1, \theta_2 \in [-\pi, \pi],$$

then for all  $\varepsilon > 0$ , we can choose our  $\delta$  as  $\delta = \varepsilon/L$ . Hence, we have  $M = \lceil L\pi/\varepsilon \rceil$  and

$$N = \left\lceil \frac{L\pi}{\varepsilon} \right\rceil \left( \left\lceil \frac{L\pi}{\varepsilon} \right\rceil + 1 \right) \left( 1 + \frac{f_{\max}}{\varepsilon} \right).$$

Next we estimate the eigenvalues of  $C_n^{-1}(f)A_n(f) - I_n$ , where  $I_n$  is the  $n$ -by- $n$  identity matrix. Using Corollary 1, Theorem 1, and the fact that

$$C_n^{-1}(f)A_n(f) - I_n = C_n^{-1}(f)(A_n(f) - C_n(f)),$$

we see that the spectrum of  $C_n^{-1}(f)A_n(f)$  is clustered around one. In particular, we have the following corollary.

**Corollary 3.** *Let  $f$  be a positive function in  $\mathcal{E}_{2\pi}$ . Then for all  $\varepsilon > 0$ , there exist  $N$  and  $M > 0$  such that for all  $n > N$ , at most  $M$  eigenvalues of the matrix  $C_n^{-1}(f)A_n(f) - I_n$  have absolute values larger than  $\varepsilon$ .*

It follows easily from the above corollary that the conjugate gradient method, when applied to the preconditioned system  $C_n^{-1}A_n$ , converges superlinearly. More precisely, for all  $\varepsilon > 0$ , there exists a constant  $c(\varepsilon) > 0$  such that the error vector  $e_q$  at the  $q$ th iteration satisfies

$$\|e_q\| \leq c(\varepsilon)\varepsilon^q\|e_0\|,$$

where  $\|x\|^2 \equiv x^*C_n^{-1/2}A_nC_n^{-1/2}x$  (see R. Chan and Strang [1] for a proof). Thus, the number of iterations to achieve a fixed accuracy remains bounded as the matrix order  $n$  is increased. Since each iteration requires  $O(n \log n)$  operations using the Fast Fourier Transform (see Strang [11]), the work of solving the equation  $A_nx = b$  to a given accuracy  $\delta$  is  $c(f, \delta)n \log n$ , where  $c(f, \delta)$  is another constant that depends only on  $f$  and  $\delta$ .

#### 4. NUMERICAL RESULTS

In this section, we test the convergence rate of the preconditioned systems with generating functions in  $\mathcal{E}_{2\pi}$ . A possible candidate is the Hardy-Littlewood series  $H(\theta)$  given by (4). However, we note that  $H(\theta)$  is not a positive function in  $[-\pi, \pi]$ . In fact, we find numerically that when  $n = 512$ , the minimum of the function

$$H_n(\theta) \equiv \sum_{k=1}^n \left\{ \frac{e^{ik \log k}}{k} e^{ik\theta} + \frac{e^{-ik \log k}}{k} e^{-ik\theta} \right\}$$

is approximately equal to  $-4.146$ . Thus, we choose the function  $H(\theta) + 4.2$  as the generating function for our numerical experiments. Three circulant preconditioners are tested, namely T. Chan's preconditioner  $C_n$  with diagonals  $c_k$  given by (2), R. Chan's preconditioner  $R_n$  with diagonals  $r_k$  given by

$$r_k = \begin{cases} a_k + a_{k-n}, & 0 \leq k < n, \\ \bar{r}_{-k}, & 0 < -k < n, \end{cases}$$

TABLE 1  
Number of iterations for different systems

$n$	$A_n$	$C_n^{-1}A_n$	$R_n^{-1}A_n$	$S_n^{-1}A_n$
16	13	8	8	8
32	18	10	10	9
64	27	11	9	9
128	43	11	9	9
256	51	10	9	9
512	58	9	9	9

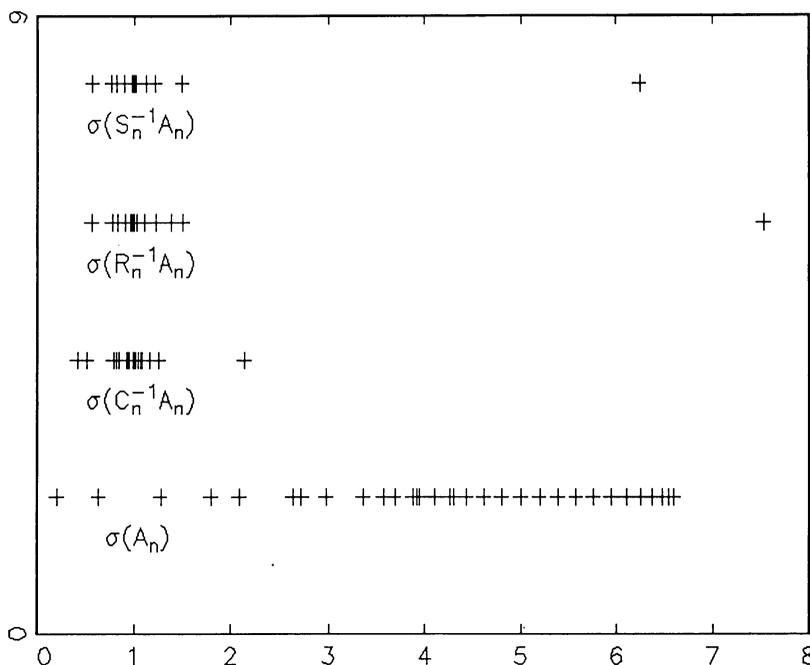


FIGURE 1  
Spectra of the preconditioned systems for  $n = 32$

and Strang's preconditioner  $S_n$  with diagonals  $s_k$  given by

$$s_k = \begin{cases} a_k, & 0 \leq k < n/2, \\ a_{k-n}, & n/2 \leq k \leq n, \\ \bar{s}_{-k}, & 0 < -k < n. \end{cases}$$

The spectra of  $A_n$ ,  $C_n^{-1}A_n$ ,  $R_n^{-1}A_n$ , and  $S_n^{-1}A_n$  for  $n = 32$  are presented in Figure 1. Table 1 shows the number of iterations required to make  $\|r_q\|_2/\|r_0\|_2 < 10^{-7}$ , where  $r_q$  is the residual vector after  $q$  iterations. The right-hand side  $b$  is the vector of all ones and the zero vector is our initial guess. The computations are done by using 8-byte arithmetic on a Vax 6420. We see that as  $n$  increases, the number of iterations increases for the original matrix  $A_n$ , while it stays almost the same for the preconditioned matrices. Moreover, all preconditioned systems converge at the same rate for large  $n$ . As for the time comparison, we report that for  $n = 512$ , it requires about 8.5 seconds to solve the original system and about 1.5 seconds to solve the preconditioned systems. Thus there is an increase in speed by a factor of about 5 to 6 when preconditioning is employed.

#### BIBLIOGRAPHY

1. R. Chan and G. Strang, *Toeplitz equations by conjugate gradients with circulant preconditioner*, SIAM J. Sci. Statist. Comput. **10** (1989), 104-119.
2. R. Chan, *The spectrum of a family of circulant preconditioned Toeplitz systems*, SIAM J. Numer. Anal. **26** (1989), 503-506.
3. —, *Circulant preconditioners for Hermitian Toeplitz systems*, SIAM J. Matrix Anal. Appl. **10** (1989), 542-550.

4. R. Chan, X. Jin, and M. Yeung, *The circulant operator in the Banach algebra of matrices*, *Linear Algebra Appl.* **149** (1991), 41–53.
5. —, *The spectra of super-optimal circulant preconditioned systems*, *SIAM J. Numer. Anal.* **28** (1991), 871–879.
6. T. Chan, *An optimal circulant preconditioner for Toeplitz systems*, *SIAM J. Sci. Statist. Comput.* **9** (1988), 766–771.
7. E. Cheney, *Introduction to approximation theory*, McGraw-Hill, New York, 1966.
8. U. Grenander and G. Szegő, *Toeplitz forms and their applications*, 2nd ed., Chelsea, New York, 1984.
9. T. Ku and C. Kuo, *Design and analysis of Toeplitz preconditioners*, *IEEE Trans. Acoust. Speech Signal Process.* (to appear).
10. W. Rudin, *Principles of mathematical analysis*, 3rd ed., McGraw-Hill, New York, 1985.
11. G. Strang, *A proposal for Toeplitz matrix calculations*, *Stud. Appl. Math.* **74** (1986), 171–176.
12. E. Tyrtyshnikov, *Optimal and super-optimal circulant preconditioners*, *SIAM J. Matrix Anal. Appl.* (to appear).
13. J. Wilkinson, *The algebraic eigenvalue problem*, Clarendon Press, Oxford, 1965.
14. A. Zygmund, *Trigonometric series*, 2nd ed., Cambridge Univ. Press, 1959.

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF HONG KONG, HONG KONG