

# A preconditioner for systems with symmetric Toeplitz blocks

S. Salapaka<sup>1†</sup>, A. Peirce<sup>2†</sup>, and M. Dahleh<sup>†</sup>

<sup>1</sup>salpax@engineering.ucsb.edu, <sup>2</sup>peirce@math.ubc.ca

<sup>†</sup>Department of Mechanical and Environmental Engineering, UCSB, Santa Barbara, CA 93106

<sup>†</sup>Department of Mathematics, The University of British Columbia, Canada

## Abstract

This paper proposes and studies the performance of a preconditioner used in the preconditioned conjugate gradient method for solving a class of symmetric positive definite systems,  $A_p x = b$ , which we call *Lower Rank Extracted Systems (LRES)*. These systems correspond to integral equations with convolution kernels defined on a union of many line segments in contrast to only one line segment in the case of Toeplitz systems. The  $p \times p$  matrix,  $A_p$ , is shown to be a principal submatrix of a larger  $N \times N$  Toeplitz matrix,  $A_N$ . The preconditioner is provided in terms of the inverse of a  $2N \times 2N$  circulant matrix constructed from the elements of  $A_N$ . The preconditioner is shown to yield clustering in the spectrum of preconditioned matrix similar to the clustering results in iterative algorithms used to solve Toeplitz systems. The analysis further demonstrates that the computational expense to solve LRE systems is reduced to  $O(N \log N)$ .

## Introduction

In this paper, we discuss the solution of a class of symmetric positive definite linear systems,  $A_p x = b$ , which we call *Lower Rank Extracted Systems (LRES)*. The coefficient matrix,  $A_p$ , has the form given by  $A_p = L_p^T A_N L_p$ , where  $A_N$  is an  $N \times N$  symmetric Toeplitz matrix and the extraction matrix,  $L_p$ , is a  $N \times p$  submatrix of an  $N \times N$  permutation matrix; i.e.,  $A_p$  is a principal submatrix of  $A_N$ . Similar to Toeplitz systems, LRES arise in the numerical modeling of convolution type integral equations. The difference is that typically in LRES, the domain of integration is a union of disjoint line segments. Therefore, Toeplitz systems, which represent the convolution type integral equations on one contiguous line segment, can be considered a special case of LRES. They appear in a wide range of scientific and engineering models, for instance in the field of image processing, in the modeling of interacting cracks, in the modeling of tabular mining excavations [1], and in the field of telecommunications in the modeling of elements in planar array antennae [2].

Their close relation to Toeplitz systems makes it possible to exploit various techniques from the vast literature for Toeplitz systems to solve them. Toeplitz systems have been studied for a long time in mathematics and engineering due

to their role in trigonometric moment problems [3], in partial differential equations, in convolution type integral equations [1], in minimum realization problems, in stochastic filtering and digital signal processing [4]. Even though most of these problems practically extend to LRES, not much attention has been given to LRES. The main contribution of this paper is that it proposes a solution to a large class of LRES which guarantees low computational expense (in the order of  $N \log N$  computations, where  $N$  is the size of the associated Toeplitz matrix  $A_N$ ).

A comprehensive survey of methods to solve Toeplitz systems (especially iterative methods) has been presented in [5]. Over the last decade, significant attention has been given to using the *Preconditioned Conjugate Gradient Method (PCGM)* [6, 7]. Many algorithms based on this method bring down the computational effort to the order of  $N \log N$  operations ([8, 9, 10, 11]). Here,  $P_N A_N \bar{x} = P_N \bar{b}$  is solved instead of  $A_N \bar{x} = \bar{b}$ . The matrix  $P_N$  is chosen so that the matrix  $P_N A_N$  has its spectrum clustered, which ensures better convergence rates.

In this paper, we use the PCGM to solve the LRES and propose a preconditioner,  $P_p$ , to solve them more efficiently. This preconditioner has been motivated by one used in [1] for solving interacting crack problems that arise in modeling mining excavations. It is remarkable that the preconditioner constructed by using the encompassing Toeplitz matrix yields such an efficient clustering of the eigenvalues associated with the multiple interacting sub-problems. In the case of Toeplitz systems, this preconditioner reduces to one of the preconditioners studied in [8]. In [8], an elegant analysis of the performance of this preconditioner for Toeplitz systems is presented. Similar preconditioners to solve Toeplitz matrices and other closely related matrices can be found in [12, 13, 14].

In section 1, we motivate the need to study LRE problems by giving an example of physical model which is represented by LRES. In section 2, we formulate the basic problem and propose the preconditioner in terms of circulant matrices. Then their properties are used to establish the clustering and convergence properties of the preconditioner for the LRES. Section 3 provides the results of some simulations. We show the persistence of the performance of the algorithm for different generating functions, different sizes of the matrices and different shapes of the domains. Finally,

in section 4 we present some concluding remarks.

### Notation

- $[T]_{p,q}$  is the element in the  $p^{th}$  row and the  $q^{th}$  column of the matrix  $T$ .
- $\delta_k$  is the kronecker-delta function,  $\delta_k = 1$  if  $k = 1$ , and  $\delta_k = 0$  if  $k \neq 0$ .
- $x^N$  is a vector of length  $2N$  given by  $(x_{-(N-1)} \dots x_0 \dots x_N)^T$ .
- $\|v\| = (\sum v_i^2)^{\frac{1}{2}}$  is the Euclidean norm of the vector  $v$ . The dimension of  $v$  is determined from the context it appears in.
- $\|T\|$  and  $\|T\|_F$  are the induced and the Frobenius (Hilbert-Schmidt) norms of the operator  $T$ .
- $T_N^x$ ,  $N \in \mathbb{N}$ ,  $x = \{x_n\}_{n=-\infty}^{\infty}$  is an  $N \times N$  symmetric Toeplitz matrix given by

$$\begin{pmatrix} x_0 & x_1 & \dots & x_{N-2} & x_{N-1} \\ x_1 & \ddots & \ddots & \ddots & x_{N-2} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & x_1 \\ x_{N-1} & \dots & \dots & x_1 & x_0 \end{pmatrix}.$$

- $H_N^x$ ,  $N \in \mathbb{N}$ ,  $x = \{x_n\}_{n=-\infty}^{\infty}$  is an  $N \times N$  symmetric Hankel matrix given by

$$\begin{pmatrix} x_1 & x_2 & \dots & x_{N-1} & x_N \\ x_2 & \ddots & \ddots & \ddots & x_{N-1} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ x_N & x_{N-1} & \dots & x_2 & x_1 \end{pmatrix}.$$

- $J_N$  is an  $N \times N$  counter identity matrix.
- $C_N^x$  is a  $2N \times 2N$  circulant matrix associated with the Toeplitz matrix  $T_N^x$  given by

$$= \begin{pmatrix} T_N^x & J_N H_N^x \\ H_N^x J_N & T_N^x \end{pmatrix}.$$

- $D_L^x(N, m, n)$  is an  $m \times n$  matrix given by  $D_L^x(N, m, n) = (I_m \ 0 \ \dots \ \dots)$

$$\begin{pmatrix} \dots & 0 & x_N & \dots & x_2 & x_1 \\ \dots & 0 & \ddots & \ddots & \ddots & x_2 \\ & & \ddots & \ddots & \ddots & \vdots \\ & & & \ddots & \ddots & x_N \\ & & & & 0 & 0 \\ & & & & \vdots & \vdots \end{pmatrix} \begin{pmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ 0 \\ I_n \end{pmatrix},$$

- $D_R^x(N, m, n)$  is an  $m \times n$  matrix given by  $D_R^x(N, m, n) = J_m D_L^x(N, m, n) J_n$ .

### 1 Motivation

In this section we give motivations for the need to study LRE systems and emphasize their relations with Toeplitz systems. For any LRE system,  $A_p x = b$ , with the domain

given by  $D = \cup_{k=1}^q V_k$ , the coefficient matrix,  $A_p$  is completely determined by the kernel of the associated integral equation and by the geometry of the domain. Accordingly, to every LRE system, we associate two matrices

1. An  $N \times N$  Toeplitz matrix,  $A_N$ , corresponding to the Toeplitz system representing an integral equation with the same kernel as the LRE system, but its domain being an interval,  $V$ , which contains the domain  $D$  of the LRE system. This matrix  $A_N$  contains all the information about the kernel.
2. An  $N \times p$  extraction matrix,  $L_p$ , that has  $q$  block-columns with the width of each column equal to the number of representative points in the corresponding segment of the domain  $D$ . In this matrix, the  $i^{th}$  block-column has only one Identity matrix (with all other entries in this block-column being 0); and every alternate block-row is a zero block. In this way, the matrix  $L_p$  has the complete information of the geometry of the domain  $D$ . Indeed, the extraction matrices,  $L_p$  are used to define the geometry of the LRE problems.

The coefficient matrix,  $A_p$  then satisfies the relation  $A_p = L_p^T A_N L_p$ .

The one dimensional integral equations and the corresponding LRE systems often represent simplified models of higher dimensional phenomena. We present one such example to emphasize the importance of the LRE systems and to understand the concepts presented above.

**Example of collinear cracks:** A simple integral equation to describe a crack located along a line (on the interval  $(a, b)$ ) in an elastic body in a state of 2D plane strain is given by

$$k \int_a^b \frac{U(\xi)}{(x - \xi)^2} d\xi = p,$$

where  $k$  is a constant depending on material properties,  $U(x)$  represents the crack opening displacement, and  $p$  represents the pressure applied to the boundary of the crack. A similar integral equation can be used to model the closure of a tabular mining excavation, whose length in the out-of-plane direction is much larger than  $b - a$  and in which the ambient stress in the rock prior to mining is given by  $-p$  (see [1] for details). A numerical approximation of this equation is obtained by partitioning the interval  $(a, b)$  into  $N$  subintervals of equal length and assuming a piecewise constant approximation to  $U(\xi)$  on each subinterval. Finding the unknowns in this discrete approximation involve solving a symmetric Toeplitz system,  $A_N x = b$ , where  $[A_N]_{i,j} = \frac{\bar{k}}{(i-j)^2 - \frac{1}{4}}$ ,  $\bar{k}$  is a constant,  $b_i = p$ , and  $x_i$  is the approximation of  $U(\xi)$  at the  $i^{th}$  element of the partition.

Similarly, an integral equation to describe  $q$  interacting collinear cracks on the intervals  $(a_1, b_1), \dots, (a_q, b_q)$  under the same physical assumptions is given by

$$k \int_{\mathcal{D}} \frac{U(\xi)}{(x - \xi)^2} d\xi = p,$$

where  $\mathcal{D}$  is the union of the intervals,  $(a_1, b_1), \dots, (a_q, b_q)$ . The numerical model for this equation, obtained by apply-

ing the same procedure as in the single crack case, yields a LRE system,  $A_p \tilde{x} = \tilde{b}$  where  $A_p$  consists of Toeplitz sub-blocks of  $A_N$ ; and  $\tilde{x}$  and  $\tilde{b}$  are subvectors of  $x$  and  $b$  respectively. In the mining context the LRE system represents the interaction of a sequence of coplanar tabular mining excavations in a state of plane strain. These coplanar mining excavations represent “rib-pillar” mining layouts commonly used in the gold mining industry.

## 2 Problem Formulation and Solution

**Problem Setting:** In the previous section, we have seen that for every LRE system, we can associate a Toeplitz matrix representative of the kernel of the integral equation. This is a many-to-one association since many LRE systems having the same kernel but different geometries can be associated with the same Toeplitz system. In this paper, we consider a sequence of Toeplitz matrices,  $\{A_N\}$ , and study the sequence,  $\{\mathcal{L}_N\}$ , of sets of LRE coefficient matrices ( $A_p$ ) that can be associated with each  $A_N$ . The sequence  $\{A_N\}$  is assumed to satisfy

- Assumptions 1** 1.  $A_N = T_N^a$  (see **Notation**), formed from the  $N$  elements  $a_0, \dots, a_{N-1}$ , of a given sequence  $\{a_n\}$  in  $\ell_1$ .
2. The sequence,  $\{a_n\}$  is such that its generating function, given by  $f(\theta) = \sum_{-\infty}^{\infty} a_k e^{ik\theta}$ , is real, symmetric, positive, and bounded away from 0; i.e.,  $\sum |a_k| < \infty$ ,  $a_k = a_{-k}$  for all  $k$  in  $\mathbb{Z}$ , and there is a  $\delta > 0$  such that  $f(\theta) > \delta > 0$  for all  $\theta$  in  $[-\pi, \pi]$ .

**Proposed Preconditioner:** We solve LRES using the Preconditioned Conjugate Gradient Method [6, 7]. In this method, a matrix (called preconditioner)  $P_p$  is designed and the system  $P_p A_p x = P_p b$  is solved instead of  $A_p x = b$ .  $P_p$  is designed so that the spectrum of  $P_p A_p$  is clustered which ensures better convergence properties (see [5]). We prescribe a preconditioner for the coefficient matrix,  $A_p$ , of an LRE system in the following way. We first form matrices  $A_N$  and  $L_p$  as in previous section and then construct a  $2N \times 2N$  circulant matrix  $C_N^a$  (see **Notation** for this construction). Since  $A_p$  is a principal submatrix of  $A_N$ , given by  $L_p^T A_N L_p$ , it is also a principal submatrix of  $C_N^a$ ; i.e.,  $A_p = \bar{L}_p^T C_N^a \bar{L}_p$ , where the *extracting matrix*,  $\bar{L}_p$  is defined by  $\bar{L}_p^T = [0 \ L_p^T]$ . Its structure is completely determined by as well as determines the geometry of the domain of the LRE system. The preconditioner,  $P_p$  is then defined by  $P_p = \bar{L}_p^T (C_N^a)^{-1} \bar{L}_p$ . In the case of Toeplitz systems,  $L_p$  is equal to  $N \times N$  Identity matrix and hence the corresponding matrix has a rank of  $N$  which is greater than any other LRE system associated with  $A_N$ . Hence the name *Lower Rank Extracted* matrices.

We have prescribed  $P_p$  in terms of the circulant matrix,  $C_N^a$ , because circulant matrices are easy to invert. Their inversion requires only  $N \log N$  multiplications and the operations can be done in parallel [15]. In a similar way the number of computations in the multiplication of a vector by  $A_p$  can be reduced (as shown in [11]). In the PCGM,

for large  $N$ , the computational effort is dominated by the preconditioner-residual product  $P_p r_j$  and the matrix vector product  $A_p d_j$  (see the PCGM algorithm in [16]). For the symmetric LRE system considered in this paper both these matrix-vector products can be easily evaluated (in  $O(N \log N)$  operations) using the corresponding circulant matrices (see [5, 1] for details).

**Clustering of the Preconditioned Matrices:** In this section, we shall show that the preconditioner that we proposed in the previous section achieves clustering of the eigenvalues. We define the clustering of the spectrum of the sequence of sets of matrices similar to the definition given for sequence of matrices in [5] by

**Definition 1** A sequence,  $\{\mathcal{L}_N\}$  of sets of matrices is said to have spectra clustered around 1 if for any given  $\epsilon > 0$ , there exist positive integers  $N_0$  and  $N_1$  such that for all  $Q_p \in \mathcal{L}_N$ ,  $N > N_0$ , at most  $N_1$  eigenvalues of the matrix  $Q_p - I_p$  have absolute value larger than  $\epsilon$ .

In this paper our aim is to show that the sequence  $\mathcal{L}_N$  of preconditioned LRE matrices ( $P_p A_p$ ) corresponding to the sequence of Toeplitz matrices  $\{A_N\}$  has spectra clustered around 1. One of the important features of our prescription is that it is given in terms of circulant matrices whose structure is exploited to attain this aim. We now present some of the important properties of these matrices (which can be easily verified) in the following proposition.

- Proposition 2.1** 1. (a) The circulant matrix  $C_N^a$  is diagonalizable, i.e.,  $C_N^a = U_N \Lambda_N U_N^T$  where  $U_N = U_N^T = U_N^{-1}$  and  $[U_N]_{ij} = \frac{1}{\sqrt{2N}} (\cos(ij\theta_N) + \sin(ij\theta_N))$  where  $0 \leq i, j \leq 2N-1$ ,  $\theta_N = \pi/N$ .
- (b)  $\Lambda_N = \text{diag}(\lambda_0^N, \dots, \lambda_{2N-1}^N)$  with  $\lambda_p^N = \sum_{k=-(N-1)}^N a_k e^{ikp\theta_N}$ ,  $0 \leq p \leq 2N-1$ .
- (c)  $\lambda_p^N = \lambda_{2N-p}^N$  for  $0 < p \leq 2N-1$ .
2. There exists an  $N_0$  in  $\mathbb{N}$  and an  $M_0$  in  $\mathbb{R}^+$  such that  $C_N^a$  is positive definite and  $\frac{1}{|\lambda_k^N|} < M_0$  for all  $N > N_0$  and  $k$  in  $\mathbb{Z}$ .
3.  $(C_N^a)^{-1} = C_N^{\xi^N}$ , where  $(\xi^N)_p =: \xi_p^N = \frac{1}{2N} \sum_{k=-(N-1)}^N \frac{1}{\lambda_k^N} e^{ipk\theta_N}$  for all  $p \in \mathbb{Z}$ .

**Relation to Fourier Coefficients of  $1/f$ :** Note that  $\lambda_j^N = \sum_{k=-(N-1)}^N a_k e^{ijk\theta_N}$  is an approximation for  $f(\theta)$  at  $j\theta_N$ ; and  $\xi_p^N = \frac{\theta_N}{2\pi} \sum_{k=-(N-1)}^N \frac{1}{\lambda_k^N} e^{ipk\theta_N}$  is a Riemann sum approximation of the integral  $\frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{1}{f(\theta)} e^{ij\theta} d\theta$ , which is the  $k^{\text{th}}$  Fourier coefficient of  $1/f$ . This suggests that the elements of  $P_p$  are approximations of the Fourier coefficients of  $g(\theta) \triangleq 1/f(\theta)$ . We establish this in the following proposition and furthermore, determine bounds on convergence rates of these approximation errors as  $N \rightarrow \infty$  under certain smoothness (and fast decay rates of the sequence  $\{a_n\}$ ) assumptions for the function,  $f$ .

- Proposition 2.2** 1. There exists a sequence  $\{\gamma_k\} \in \ell_1(-\infty, \infty)$  with  $\gamma_k = \gamma_{-k}$  such that  $g(\theta) = \sum_{k=-\infty}^{\infty} \gamma_k e^{ik\theta}$  for all  $\theta$  in  $[-\pi, \pi]$ .
2.  $\lim_{N \rightarrow \infty} N \|\gamma^N - \xi^N\|^2 = 0$  if  $\sum_{k=-(N-1)}^N |k^2 a_k^2| < \infty$ .

See [16] for the proof.

### Clustering of the Spectrum of LRE Matrices:

In this section we define a class of LRE systems and show the clustering properties of the corresponding preconditioned LRE matrices. We define a sequence of sets of LRE matrices associated with the sequence of Toeplitz matrices  $\{A_N\}$  in the following way. For every  $\epsilon > 0$ , let  $N_0(\epsilon)$  and  $N_1(\epsilon)$  be such that  $\sum_{N_0}^{\infty} k a_k^2 \leq \epsilon$ ,  $\sum_{N_0}^{\infty} k \gamma_k^2 \leq \epsilon$  and  $N \|\gamma^N - \xi^N\|^2 \leq \epsilon$  for all  $N > N_1(\epsilon)$  (this is possible by Proposition 2). Then, to every  $A_N$  for  $N > N_1(\epsilon)$ , we denote a set of LRE matrices by  $\mathcal{L}_N^A(\epsilon)$  whose elements have the form given by  $A_p = L_p^T A_N L_p$ , where

1.  $L_p$  has the structure given by

$$L_p = \begin{pmatrix} I_{p_0} & 0 & 0 & \dots \\ 0 & 0 & 0 & \dots \\ 0 & I_{p_1} & 0 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix} \begin{matrix} \} r_0 \\ \} r_1 \\ \} r_2 \\ \vdots \end{matrix},$$

and has  $n_p$  block-columns (with  $\sum_{i=0}^{n_p-1} p_i = p$ ).

2.  $r_i > N_0(\epsilon)$  for all  $0 \leq i \leq n_p - 1$ .

To each  $A_p = L_p^T T_N^A L_p \in \mathcal{L}_N^A(\epsilon)$ , we associate a preconditioner as described earlier in this paper; i.e.,  $P_p = L_p^T T_N^{\xi^N} L_p$ . The sequence of sets of preconditioned LRE matrices can be now defined by

$$\mathcal{L}_N(\epsilon) = \left\{ P_p A_p \text{ such that } A_p = L_p^T T_N^A L_p \in \mathcal{L}_N^A(\epsilon) \text{ and } P_p = L_p^T T_N^{\xi^N} L_p \right\}$$

These definitions being given, we present the following proposition,

**Proposition 2.3** If  $f \sum_{k=-\infty}^{\infty} |k^2 a_k^2| < \infty$ , and under the Assumptions 1,

for every  $\epsilon > 0$ , there exist  $N_0$  and  $N_1$  in  $\mathbb{N}$  such that

$$\|I - P_p A_p - D_p\|_F \leq \epsilon \text{ for all } P_p A_p \in \mathcal{L}_N(\epsilon) \text{ and } N \geq N_1,$$

where  $D_p$  is a block diagonal matrix which has at most  $2n_p$  non-zero  $N_0 \times N_0$  blocks.

**Proof of Proposition 2.3:** Let  $\epsilon > 0$  and  $N_0(\epsilon)$  and  $N_1(\epsilon)$  be such that  $\sum_{N_0}^{\infty} k a_k^2 \leq \epsilon$ ,  $\sum_{N_0}^{\infty} k \gamma_k^2 \leq \epsilon$  and  $N \|\gamma^N - \xi^N\|^2 \leq \epsilon$  for all  $N > N_1$  (this is possible by Proposition 2). Let  $A_p \in \mathcal{L}_N^A(\epsilon)$  for some  $N > N_1$  and  $P_p = L_p^T T_N^{\xi^N} L_p$  be its preconditioner. Note that from Proposition A.1, we have that  $P_N A_N = T_N^{\xi^N} T_N^A = I + \bar{D}$  where  $\bar{D}$  has at most 2 nonzero  $N_0 \times N_0$  blocks.

Also  $L_p^T T_N^{\xi^N} T_N^A L_p - \underbrace{L_p^T T_N^{\xi^N} L_p}_{P_p} \underbrace{L_p^T T_N^A L_p}_{A_p}$  can be rewritten as

$L_p^T T_N^{\xi^N} \tilde{L}_p \tilde{L}_p^T T_N^A L_p$  where  $\tilde{L}_p$  is such that  $L_p L_p^T + \tilde{L}_p \tilde{L}_p^T = I$ . This implies that

$$I + L_p^T \bar{D} L_p - P_p A_p = L_p^T T_N^{\xi^N} \tilde{L}_p \tilde{L}_p^T T_N^A L_p. \quad (1)$$

We first prove the following properties of  $L_p^T T_N^{\xi^N} \tilde{L}_p$  and  $\tilde{L}_p^T T_N^A L_p$  which we shall use to study the spectrum of  $P_p A_p$ ,

1.  $L_p^T T_N^{\xi^N} \tilde{L}_p = L_p^T D_N^{\gamma} \tilde{L}_p + E_N^{\gamma}$  where  $D_N^{\gamma}$  is a block tridiagonal matrix which has the form

$$\begin{pmatrix} 0 & R_{01}^{\gamma} & 0 & \dots & \dots & 0 \\ L_{10}^{\gamma} & 0 & R_{12}^{\gamma} & \ddots & \ddots & \vdots \\ 0 & \ddots & 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & R_{n_r-2, n_r-1}^{\gamma} \\ 0 & \dots & \dots & 0 & L_{n_r-1, n_r-2}^{\gamma} & 0 \end{pmatrix},$$

where  $R_{ij}^{\gamma} = D_R^{\gamma}(N_0, r_i, r_j)$  and  $L_{ij}^{\gamma} = D_L^{\gamma}(N_0, r_i, r_j)$  (see Notation); and  $\|E_N^{\gamma}\|_F \leq 2n_p(n_r - n_p)\epsilon$ .

2.  $\tilde{L}_p^T T_N^A L_p = \tilde{L}_p^T D_N^A L_p + E_N^A$  where  $\|E_N^A\|_F \leq 2n_p(n_r - n_p)\epsilon$  and  $D_N^A$  is defined in the same way as  $D_N^{\gamma}$ .
3.  $D_N^{\gamma} D_N^A$  is a block diagonal matrix with only  $2n_r - 2$  non-zero  $N_0 \times N_0$  blocks.

1), 2) Consider the product  $L_p^T T_N^{\xi^N} \tilde{L}_p$ . It is independent of the  $n_r$  diagonal blocks ( $r_i \times r_i$  blocks,  $0 \leq i \leq n_r - 1$ ) in  $T_N^{\xi^N}$  due to orthogonality of the matrices  $L_p$  and  $\tilde{L}_p$ .

Therefore these diagonal blocks in  $T_N^{\xi^N}$  can be replaced by zeros and still the product remains unchanged. Therefore this product can be rewritten as  $L_p^T \bar{D}_N^{\xi^N} \tilde{L}_p$ , where  $\bar{D}_N^{\xi^N}$  is obtained by substituting the diagonal blocks (the  $n_r \times r_i$  blocks) in  $T_N^{\xi^N}$  by zero blocks. Also note that  $\tilde{R}_i^{\xi^N}$  and  $\tilde{L}_i^{\xi^N}$  are submatrices of  $J_N H_N^{\xi^N}$  and  $H_N^{\xi^N} J_N$  respectively. From our choice of  $N_1$ , we have  $\|J_N H_N^{\xi^N} - J_N H_N^{\gamma}\|_F \leq \epsilon$  which implies  $\|\tilde{R}_i^{\xi^N} - \tilde{R}_i^{\gamma}\|_F \leq \epsilon$  (using Lemmas 1 and 2 in [16]). Therefore

$$\|\tilde{R}_i^{\xi^N} - (R_{ij}^{\gamma} \ 0)\|_F \leq \|\tilde{R}_i^{\xi^N} - \tilde{R}_i^{\gamma}\|_F + \|\tilde{R}_i^{\gamma} - (R_{ij}^{\gamma} \ 0)\|_F \leq 2\epsilon,$$

where the zero block is of appropriate size. Similarly, we have  $\|\tilde{L}_i^{\xi^N} - (0 \ L_{ij}^{\gamma})\|_F \leq 2\epsilon$ . Therefore

$$\begin{aligned} \underbrace{\|L_p^T T_N^{\xi^N} \tilde{L}_p - L_p^T D_N^{\gamma} \tilde{L}_p\|_F}_{\triangleq E_N^{\gamma}} &= \|L_p^T (\bar{D}_N^{\xi^N} - D_N^{\gamma}) \tilde{L}_p\|_F \\ &\leq 2n_p(n_r - n_p)\epsilon. \end{aligned}$$

- 2) This can be proved in the same way as (1).

3) Note that the matrices  $R_{ij}^{\gamma}$  and  $R_{ij}^A$  have only a lower left non-zero  $N_0 \times N_0$  block; and the matrices  $L_{ij}^{\gamma}$  and  $L_{ij}^A$  have only a top right non-zero  $N_0 \times N_0$  block. This structure on them implies that the products  $R_{ij}^{\gamma} R_{i'j'}^A = 0$  and  $L_{ij}^{\gamma} L_{i'j'}^A = 0$ ; and the products  $R_{ij}^{\gamma} L_{i'j'}^A$  and  $L_{ij}^{\gamma} R_{i'j'}^A$  are block diagonal with only one non-zero  $N_0 \times N_0$  block for all  $0 \leq i, j, i', j' \leq$

$n_r - 1$ . This implies that the product  $D_N^\gamma D_N^\alpha$  is a block diagonal matrix ( $= \text{diag}(R_{01}^\gamma L_{10}^\alpha, L_{10}^\gamma R_{01}^\alpha + R_{12}^\gamma L_{21}^\alpha, \dots)$ ) with at most  $2n_r - 2$  non-zero  $N_0 \times N_0$  blocks. Therefore  $L_p^\gamma D_N^\gamma D_N^\alpha L_p$  is a block diagonal matrix with at most  $2n_p - 2$  non-zero  $N_0 \times N_0$  blocks.

Also, note that from the structures of  $D_N^\alpha$  and  $D_N^\gamma$ , we have  $L_p^\gamma D_N^\gamma L_p = 0$  and  $L_p^\alpha D_N^\alpha L_p = 0$  and therefore  $L_p^\gamma D_N^\gamma L_p L_p^\alpha D_N^\alpha L_p = 0$ . Now Equation 1 can be further simplified as

$$\begin{aligned} I + L_p^\gamma \tilde{D} L_p - P_p A_p &= (L_p^\gamma D_N^\gamma \tilde{L}_p + E_N^\gamma)(\tilde{L}_p^\alpha D_N^\alpha L_p + E_N^\alpha) \\ &= L_p^\gamma \underbrace{D_N^\gamma D_N^\alpha L_p}_{\triangleq \tilde{D}} + E_N^\gamma(\tilde{L}_p^\alpha D_N^\alpha L_p + E_N^\alpha) + (L_p^\gamma D_N^\gamma \tilde{L}_p)E_N^\alpha. \end{aligned}$$

Therefore,

$$\|I - P_p A_p - \underbrace{L_p^\gamma(\tilde{D} + \tilde{D})L_p}_\triangleq D_p\|_F \leq M\epsilon,$$

where  $M$  is an upper bound on  $2n_p(n_r - n_p)(\|C_N^\alpha\|_2 + \|(C_N^\alpha)^{-1}\|_2)$ , and  $D_p$  is a matrix with at most  $2n_p$  non-zero  $N_0 \times N_0$  blocks. As  $\epsilon > 0$  and  $N > N_1$  were chosen arbitrarily, we have proved the proposition. ■

In the following proposition, we use the positive definiteness of  $A_p$  and Theorem 2 in [17] to prove the following Proposition,

**Proposition 2.4** *If  $\sum_{k=-\infty}^{\infty} |k^2 a_k^2| < \infty$ , and under the Assumptions 1,*

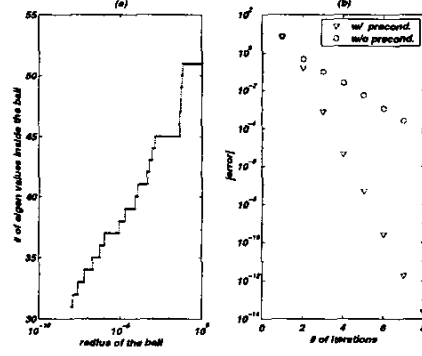
*for every  $\epsilon > 0$ , there exist  $N_0$  and  $N_1$  in  $\mathbb{N}$  such that there are at least  $p - 2n_p N_0$  eigenvalues  $\alpha_j$ , of  $P_p A_p$  satisfying  $|\alpha_j - 1| \leq 4\epsilon^2$  for all  $P_p A_p \in \mathcal{L}_N(\epsilon)$  and  $N > N_1$ ; i.e.,*

*the spectrum of the sequence of sets of preconditioned LRE matrices,  $\{\mathcal{L}_N(\epsilon)\}$  clusters around 1 for every  $\epsilon > 0$ . See [16] for the proof.*

### 3 Simulation Results

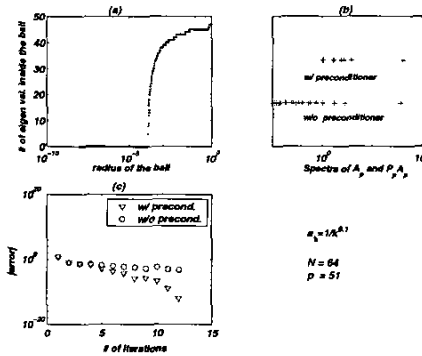
We have shown in previous sections that the preconditioners,  $P_p$ , which are extracted from the inverse of the circulant matrix  $C_N^\alpha$  yield spectra of  $\{P_p A_p\}$  that are clustered around 1. This is desirable from a computational point of view as circulant matrices are easy to invert in addition to the fact that the PCGM converges more rapidly if the eigenvalues are clustered (see [7]).

In Figures 1-2, we provide simulation results for LRE systems with different kernels but defined on the same domain. Each of these LRE systems,  $A_p x = b$ , is described by a  $51 \times 51$  coefficient matrix,  $A_p$  which is a principal submatrix of a corresponding  $64 \times 64$  Toeplitz matrix,  $A$ ; i.e.  $A_p = L_p^T A L_p$ . The domain consists of 3 line segments and the structure of  $L_p$  (as laid out in Section 2) is completely specified by the dimensions,  $r_0 = 17$ ,  $r_1 = 7$ ,  $r_2 = 17$ ,  $r_3 = 6$  and  $r_4 = 21$ . The kernels of the integral equations are specified by a different  $A$  in each of these LRE systems. In all simulations, we assumed that the given vector,  $b = [1 \ 1 \ \dots]^T$ .



**Figure 1:** The three cracks problem ( $N = 64$ ,  $p = 51$ ): (a) Plot of number of eigenvalues (of  $P_p A_p$ ) within a ball around 1 versus the radius of the ball. (b) Comparison of the convergence rates of PCGM between the preconditioned and non-preconditioned cases.

In Figure 1, the kernel,  $A$  is generated by the function  $f(\theta) = 2\pi \sin(|\theta|/2)$  whose fourier coefficients form the sequence  $\{\frac{-1}{n^2-1}\}$ . This problem represents the interaction of three cracks (see the mining example in Section 1). In (a), we observe that a majority (40 out of 51) of the eigenvalues are clustered around 1 (within a radius of  $10^{-4}$ ). This clustering of the eigenvalues is exploited by the PCGM and the rapid convergence of the PCGM can be observed in (b). These convergence trends were found to persist on simulations performed with other generating functions (such as  $f(\theta) = \theta^4 + 1$ ) and sequences (e.g.  $\{(-1)^{\text{rand}(k)}/k^2\}$ , where  $\text{rand}(k) \in \mathbb{N}$  is obtained by truncating a random real number,  $r$  ( $0 < r < 1000$ ) generated by using the “rand” function in MATLAB.)



**Figure 2:** Robustness of the PCGM ( $N = 64$ ,  $p = 51$ ): (a) Plot of number of eigenvalues (of  $P_p A_p$ ) within a ball around 1 versus the radius of the ball. (b) Clustering of eigenvalues of  $P_p A_p$  compared to that of  $A_p$ . (c) Comparison of the convergence rates of PCGM between the preconditioned and non-preconditioned cases.

Figure 2 shows results for a LRE system generated by a se-

quence  $a_k = k^{-0.1}$ . Note that the corresponding series is not even convergent, but still the algorithm works satisfactorily. The fact that these algorithms continue to perform well even beyond the limitations of our analysis demonstrates their robustness.

The clustering, convergence and robustness properties were found to persist in many other simulations (which we do not present here) that we did by changing the kernels, domain geometries and matrix sizes. For example, similar trends were found in the simulation of the interaction of six cracks problem with same generating function as in the three crack problem but defined on a completely different domain. Here  $A_p$  was chosen to be a  $90 \times 90$  matrix and the corresponding  $A$  to be a  $128 \times 128$  matrix.

#### 4 Conclusions

In this paper we have introduced and analyzed preconditioners ( $P_p$ ) in PCGM for the efficient solution of Lower Rank Extracted systems (LRES),  $A_p x = b$ . The elements of the preconditioners are shown to approximate the Fourier coefficients of the reciprocal of the generating function associated with the LRE system. Under fairly mild assumptions on the generating function,  $f(\theta)$  or alternatively on the generating sequence  $\{a_N\}$  these properties are exploited in order to prove clustering of the eigenvalues of the matrices  $P_p A_p$ . Also, these systems are shown to be subsystems of Toeplitz systems,  $A_N x = b$ . For LRES, the PCGM converges to a specified tolerance in  $O(N \log N)$  operations where  $N$  is the size of  $A_N$ . To study the preconditioner,  $P_p$ , many simulation of LRES with different kernels, sizes and domains have been presented. Simulation results corroborate the theoretical findings regarding clustering of the spectra of preconditioned matrices and the associated convergence rates. In particular, the majority of the eigenvalues of  $P_p A_p$  fall in the vicinity of 1. In addition, the simulations demonstrate that the algorithm is robust in that it still yields significant clustering even for Toeplitz matrices derived from sequences which did not satisfy the restrictions imposed by the hypotheses of the propositions presented. This indicates that theoretical results established in this paper might be proved under more relaxed conditions.

#### A Appendix

Here, we present a proposition that shows that in the case of Toeplitz matrices, the product  $I - P_N A_N$  can be approximated by a block diagonal matrix with a large 0 block. This case has been analyzed [8] and the proposition presented here is very similar to (and can be derived from) Lemma 7 in [8].

**Proposition A.1** *If  $\sum_{k=-\infty}^{\infty} |k^2 a_k^2| < \infty$ , then for  $\epsilon > 0$  there exist  $N_0$  and  $N_1$  in  $\mathbb{N}$  such that*

$$\|H_N^{\epsilon N} H_N^a - D_N\|_F \leq \epsilon \text{ for all } N \geq N_1.$$

*where  $D_N$  is a block diagonal matrix with only two non-zero  $N_0 \times N_0$  blocks.*

#### References

- [1] A. P. Peirce, S. Spottiswoode, and J.A.L. Napier. The Spectral Boundary Element Method: A New Window on Boundary Elements in Rock Mechanics. *Int. J. Rock Mech. Min. Sci. and Geomech. Abstr.*, 29(4):379–400, 1992.
- [2] Zhuang Y. On the use of hybrid full-wave cg-fft and spectrally preconditioned cg-fft methods for analyzing large microstrip antenna arrays. *Ph. D. Thesis, Mc Master University*, 7, August 1995.
- [3] N.I. Akhiezer. *The Classical Moment Problem*. Oliver and Boyd, 1st edition, 1965.
- [4] T. Kailath. A View of Three Decades of Linear Filtering Theory. *IEEE Trans. Information Theory*, IT-20:145–181, 1974.
- [5] R.H. Chan and M.K. Ng. Conjugate Gradient Methods for Toeplitz Systems. *SIAM Review*, 38(3):427–482, September 1996.
- [6] A. Quarteroni, R. Sacco, and S. Fausto. *Numerical Mathematics*. Springer, 1st edition, 2000.
- [7] D.G. Luenberger. *Linear and Nonlinear Programming*. Addison-Wesley Publishing Company, 2nd edition, 1984.
- [8] R.H. Chan and K.P. Ng. Toeplitz Preconditioners for Hermitian Toeplitz Systems. *Linear Algebra and its Applications*, 190:181–208, September 1993.
- [9] G Strang. A Proposal for Toeplitz Matrix Calculations. *Stud. Appl. Math.*, 74:171–176, 1986.
- [10] T Chan. An Optimal Circulant Preconditioner for Toeplitz Systems. *SIAM J. Numer. Anal.*, 30:1193–1207, 1993.
- [11] R. Chan. Circulant Preconditioners for Hermitian Toeplitz Systems. *SIAM J. Matrix Anal. Appl.*, 10:542–550, 1989.
- [12] M. Hanke and J. Nagy. Toeplitz-approximate inverse preconditioner for banded toeplitz matrices. *Numerical Algorithms*, 7:183–199, 1994.
- [13] J.G. Nagy, R.J. Plemmons, and T.C. Torgersen. Iterative Image Restoration using Approximate Inverse Preconditioning. *IEEE Transactions on Image Processing*, 5(7):1151–1161, July 1996.
- [14] T. Chan and J. Olkin. Circulant preconditioners for toeplitz-block matrices. *Numerical Algorithms*, 6:89–101, 1994.
- [15] A.V. Oppenheim and R.W. Schaffer. *Discrete-Time Signal Processing*. Prentice Hall, Englewood Cliffs, New Jersey, 1st edition, 1989.
- [16] S. Salapaka, A. Peirce, and M. Dahleh. A preconditioner for systems with symmetric Toeplitz blocks. Preprint at <http://www.engr.ucsb.edu/~salpaz/cdc2001/preprint.pdf>.
- [17] W. Kahan. Spectra of nearly Hermitian Matrices. *Proceedings of American Mathematical Society*, 48(1):11–17, March 1975.