TOEPLITZ EQUATIONS BY CONJUGATE GRADIENTS WITH CIRCULANT PRECONDITIONER*

RAYMOND H. CHAN† AND GILBERT STRANG‡

Abstract. This paper studies the solution of symmetric positive definite Toeplitz systems Ax = b by the preconditioned conjugate gradient method. The preconditioner is a circulant matrix C that copies the middle diagonals of A, and each iteration uses the Fast Fourier Transform. Convergence is governed by the eigenvalues of $C^{-1}A$ —a Toeplitz-circulant eigenvalue problem—and it is fast if those eigenvalues are clustered. The limiting behavior of the eigenvalues is found as the dimension increases, and it is proved that they cluster around $\lambda = 1$. For a wide class of problems the error after q conjugate gradient steps decreases as r^{q^2} .

Key words. Toeplitz, circulant, conjugate gradient, Hankel

AMS(MOS) subject classification. 65F

1. Introduction. In this paper we discuss a class of linear systems Ax = b. The matrix A has the **Toeplitz Property:** Down each diagonal its entries are constant. The *i*, *j* entry is a_{i-j} , and we assume symmetry and positive definiteness. Such systems are fundamental in signal processing and time series, where the convolution form reflects invariance in time or in space (stationarity or homogeneity). Toeplitz matrices also arise directly from constant-coefficient partial differential equations, and from integral equations with a convolution kernel, when those equations are made discrete.

With periodicity, these problems can be solved quickly by Fourier transform. The convolution becomes a multiplication and deconvolution is straightforward. In the nonperiodic case, which is analogous to a problem on a finite interval (or on a bounded region in the multidimensional case), this direct solution is lost. The inverse of a Toeplitz matrix is not Toeplitz, because of the presence of a boundary and the absence of periodicity. Nevertheless the matrix A is determined by only n coefficients a_0, \dots, a_{n-1} , rather than by n^2 . Algorithms that exploit the Toeplitz property are much faster than the $n^3/6$ operations of symmetric elimination, and direct methods based on the Levinson recursion formula [1] are in constant use. A number of superfast methods have been created in the last ten years, and an implementation by Ammar and Gragg [2] is giving excellent results. (See [2]-[5] for references, and note also the algorithms developed for systolic arrays [6].) More recently, the second author proposed an iterative method [7] that, it is hoped, will be fast and flexible. We report at the end on recent experiments, after an analysis of this iterative method.

Note. The Levinson algorithm conventionally uses $2n^2$ multiplications, but the "split-Levinson" form [8] reduces the count to nearly n^2 . The superfast methods are $O(n \log^2 n)$, and the constant in the new implementation based on Schur's algorithm is about 8. It has become competitive with Levinson at reasonable n—a remarkable achievement. Our iterative algorithm needs $O(n \log n)$ operations per step from the **Fast Fourier Transform** (FFT), and the number of steps depends (inevitably) on the

^{*} Received by the editors May 29, 1987; accepted for publication (in revised form) March 16, 1988.

[†] Department of Mathematics, University of Hong Kong, Hong Kong. The research of this author was supported by National Science Foundation grants DCR84-05506 and DCR86-02563.

[‡] Department of Mathematics, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139. The research of this author was supported by Army Research Office grants DAAG29-83-K0025 and DAAL03-86-K0171, and National Science Foundation grant DMS87-03313.

Toeplitz matrix and the accuracy required. The goal of this paper is to analyze the rate of convergence in terms of the function $f = \sum a_k z^k$ constructed from the matrix, and to show that for analytic f (and many other functions) the number of iterations is independent of n. Tests on serious applications (to time series, scattering, or exploration data) are still in the future.

The iterative method uses a preconditioner. The Toeplitz matrix is replaced by a *circulant matrix*. It retains the Toeplitz property and adds periodicity. Each diagonal in the lower triangular part wraps around into a diagonal in the upper triangular part, and the entries satisfy $c_{ij} = c_{i-j} = c_{i-j+n}$. The distinction between Toeplitz and circulant matrices is seen (in the symmetric case) in

$$A = \begin{bmatrix} a_0 & a_1 & \cdot & a_{n-2} & a_{n-1} \\ a_1 & a_0 & a_1 & \cdot & a_{n-2} \\ a_2 & a_1 & a_0 & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & a_1 \\ a_{n-1} & \cdot & a_2 & a_1 & a_0 \end{bmatrix} \text{ and } C = \begin{bmatrix} c_0 & c_1 & \cdot & c_2 & c_1 \\ c_1 & c_0 & c_1 & \cdot & c_2 \\ c_2 & c_1 & c_0 & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & c_1 \\ c_1 & \cdot & c_2 & c_1 & c_0 \end{bmatrix}$$

The diagonals containing c_1 reappear in the corners, where the matrix A has a new (and probably smaller) entry a_{n-1} . To go from A to C will require changing about $n^2/4$ entries, and the key question in analyzing convergence will be the eigenvalues of $C^{-1}A$.

Multiplication by a circulant C is identical to discrete convolution. The linear system Cz = b is the convolution equation c * z = b, where c is the first column of C. After a discrete Fourier transform it becomes $\hat{c}\hat{z} = \hat{b}$. Therefore \hat{z} is given by a component-by-component division, and z is recovered from the inverse Fourier transform. The components of \hat{c} are proportional to the eigenvalues of C, and this convolution rule is the diagonalization of the circulant matrix [9]: $z = F^{-1}\Lambda^{-1}Fb$. This is one of the rare instances in which a linear system is solved by diagonalizing the coefficient matrix! Normally elimination is much faster, but the Fourier matrix F (its entries are the complex roots of unity $F_{jk} = w^{jk} = \exp 2\pi i jk/n$, and its columns are the eigenvectors of every circulant matrix) is very special.

The speed of the iterative method depends on the fact that multiplication by F and F^{-1} —the discrete Fourier transform and its inverse—can be done so quickly. Those multiplications are computed by the FFT, which dominates each step of the iteration. It requires only $n \log n$ multiplications, and the calculations can be done in parallel. It applies directly to C and our goal is to apply it also to Ax = b—reaching the required tolerance in a number of steps which in the best case is independent of n.

We mention that *multiplication* by a Toeplitz matrix A (but not inversion) is also quick by the FFT. The matrix is extended to a circulant A^* of order 2n, the vector d is completed to d^* by n zeros, and Ad appears in the first n components of A^*d^* , which is another discrete convolution. The goal is to replace A by C in any linear system to be solved, and to use A itself only in matrix multiplications.

This is exactly what is achieved by the ordinary iterative method $Cx_{n+1} = (C-A)x_n + b$, and also by the **preconditioned conjugate gradient method.** There is a Toeplitz multiplication on the right side and a circulant inversion on the left. We will see that the ordinary iterations can diverge; they depend on the extreme eigenvalues of $C^{-1}A$, which are not in close control. However the conjugate gradient method can be very effective. Its convergence rate also depends on the eigenvalues λ_i of $C^{-1}A$, but not exclusively on λ_1 and λ_n . Conjugate gradient convergence is fast when the eigenvalues are clustered, and that is the property established in this paper. Thus we

want to show that the circulant matrix satisfies, for large n, the following two essential requirements for a good preconditioner:

- (1) Cz = d can be solved quickly and stably (Theorem 1 will estimate $||C^{-1}||$).
- (2) C is close to A (the eigenvalues of $C^{-1}A$ are clustered near 1).

For completeness we list the steps of the preconditioned conjugate gradient method, which gives the exact solution at step n; however, it is treated as an iterative method and stops earlier. Each iteration contains the periodic linear system with coefficient matrix C, the multiplication by A, and the two inner products that appropriately orthogonalize the directions d_j . Starting from $x_0 = 0$ and $r_0 = b$,

Solve $Cz_{j-1} = r_{j-1}$,

 $\beta_{j} = z_{j-1}^{T} r_{j-1} / z_{j-2}^{T} r_{j-2} \quad (\text{except } \beta_{1} = 0),$ $d_{j} = z_{j-1} + \beta_{j} d_{j-1} \qquad (\text{except } d_{1} = z_{0}),$ $\alpha_{j} = z_{j-1}^{T} r_{j-1} / d_{j}^{T} A d_{j},$ $x_{j} = x_{j-1} + \alpha_{j} d_{j},$ $r_{j} = r_{j-1} - \alpha_{j} A d_{j}.$

2. The eigenvalues of $C^{-1}A$. We want to choose C close to A. The simplest construction is to copy the central diagonals of A and bring them around to complete the circulant. Starting from the first column a_0, \dots, a_{n-1} of A, with n = 2m, the first column of C is $a_0, \dots, a_m, \dots, a_1$. If A decays quickly away from the main diagonal, then C starts to do the same but increases again as we approach the corner. By substituting the vector $(1, 0, \dots, 0, -1)$ into the Rayleigh quotient x^TAx/x^TCx , we see that the largest eigenvalue of $C^{-1}A$ is at least

(1)
$$\frac{a_0 - a_{n-1}}{a_0 - a_1} \leq \lambda_{\max} (C^{-1} A).$$

This can easily exceed 2, in which case the ordinary iteration $Cx_{n+1} = (C - A)x_n + b$ will fail. The iterating matrix $I - C^{-1}A$ has $1 - \lambda_{max}$ outside the unit circle. However, the conjugate gradient method can compensate for any single outlying eigenvalue in a single iteration. The question is whether many other eigenvalues are far from unity, when corners of order m = n/2 are different in C and A.

Our first results are experimental [7]. With diagonal entries $a_k = 1/(1+k)$ the eigenvalues for n = 12 are .707, .957, \cdots , 1.047, 1.880. The largest and smallest make ordinary iteration too slow, but the other eigenvalues are clustered around 1. As the order n is increased, they approach limiting values. It is not always clear numerically, say for λ_4 , whether the limit is 1. In these experiments, and in others with different diagonals a_k , the x_i converge quickly to $x = A^{-1}b$.

The next results are theoretical [10]. The matrices with geometrically decreasing diagonals $a_k = t^k$ exhibit a remarkable pattern in the computations, and the eigenvalues (and eigenvectors) of $C^{-1}A$ can be verified analytically. The extremes are 1/(1+t) and 1/(1-t), and $\lambda = 1$ is a double eigenvalue. What is striking is that *there are only two other eigenvalues of* $C^{-1}A$. For a matrix of order 1024, each is repeated 510 times! Those eigenvalues are exponentially close to 1: $\lambda = 1/(1+t^{n/2})$ and $\lambda = 1/(1-t^{n/2})$. Convergence of conjugate gradients is extremely fast, and our goal is to see when this exponential clustering can be predicted.

These matrices A were studied by Kac, Murdock, and Szegö [11], who observed that they have a special property: A^{-1} is tridiagonal. The generating function of A is

(2)
$$f = \sum_{-\infty}^{\infty} a_k e^{ik\theta} = \sum_{-\infty}^{\infty} t^{|k|} e^{ik\theta} = \frac{1-t^2}{(1-t e^{i\theta})(1-t e^{-i\theta})}.$$

It is real and positive, so that A is symmetric positive definite (for |t| < 1). It is the reciprocal of a three-term polynomial, which underlies the fact that A^{-1} is banded and that only two limiting eigenvalues are different from 1.

The ideal approach is to learn about the spectrum of $C^{-1}A$ from this function $f(\theta)$. We recognize that in practice the matrices are finite and the very distant diagonals a_k will not be used. But the asymptotic properties appear to be decisive, and all the information about C and A is in f. The goal is to turn a problem in operator theory into a problem in function theory.

Remark. The eigenvalues of the Toeplitz matrix alone have been studied in detail [12]-[14] and their behavior is entirely different. Instead of clustering around 1, they are "equidistributed" with the values of f itself [15]. The same is true for the circulants alone [16]. The new problem is to study the product $C^{-1}A$, and this paper computes the limits of the eigenvalues as $n \to \infty$. We believe it will also be a key to the finite case—where we look first at examples.

3. Uniform invertibility of C. Suppose that the Toeplitz matrices A_n , of order n, are finite sections of a fixed singly infinite positive definite matrix A_{∞} . The *i*, *j* entries of A_n and A_{∞} are a_{i-j} , and the associated function

$$f(\theta) = \sum_{-\infty}^{\infty} a_k e^{ik\theta}$$

is real and positive. We will assume that the sequence a_k is in l^1 , so that f belongs to the Wiener class: $\sum |a_k| < \infty$. Then the function 1/f associated with A_{∞}^{-1} belongs to the same class, and a more precise analysis becomes possible.

The first step is to consider C. The construction that copies the middle diagonals of A does not guarantee the invertibility of C.

Example 1.

$$A = \begin{bmatrix} 2 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 2 \end{bmatrix} \text{ and } C = \begin{bmatrix} 2 & -1 & 0 & -1 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ -1 & 0 & -1 & 2 \end{bmatrix}.$$

For this "second-difference matrix" the construction produces a singular C. That occurs whenever A comes from discretizing an operator with no zero-order term. Dirichlet boundary conditions leave A invertible, while periodic boundary conditions make C singular. It is a case when sines should replace exponentials. The matrix S with entries $\sin jk\pi/(n+1)$ has the eigenvectors of this A as its columns. The Fast Sine Transform carries out multiplication by S and S^{-1} in n log n steps. Therefore the new preconditioner can be SDS^{-1} , where the diagonal matrix D has entries $d_j = \sum a_k e^{ijk\theta}$, $\theta = \pi/(n+1)$.

In this example A_{∞} is only semidefinite, and $f(\theta) = 2 - 2 \cos \theta$ is only nonnegative. Iteration is not needed because the matrix is banded, but for a full matrix with $\sum a_k = 0$ the idea may be useful—conjugate gradients preconditioned by a sine transform. In this paper we stay with the positive definite case f > 0.

Example 2.

$$A = \begin{bmatrix} .7 & \frac{1}{2} & \frac{1}{4} & \frac{1}{8} \\ \frac{1}{2} & .7 & \frac{1}{2} & \frac{1}{4} \\ \frac{1}{4} & \frac{1}{2} & .7 & \frac{1}{2} \\ \frac{1}{8} & \frac{1}{4} & \frac{1}{2} & .7 \end{bmatrix} \text{ and } C = \begin{bmatrix} .7 & \frac{1}{2} & \frac{1}{4} & \frac{1}{2} \\ \frac{1}{2} & .7 & \frac{1}{2} & \frac{1}{4} \\ \frac{1}{4} & \frac{1}{2} & .7 & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{2} & .7 \end{bmatrix}.$$

The smallest eigenvalue of A is 3/40, whereas C has the eigenvalue -1/20. The new feature is that this A extends to a positive definite A_{∞} —it is the Kac-Murdock-Szegö matrix with diagonals $a_k = (\frac{1}{2})^k$ and with a_0 changed to .7. Note that our algorithm will recognize the indefiniteness of C at the first step, when $Cz_0 = r_0$ is solved by the FFT. C is diagonalized so its eigenvalues are made explicit, and the algorithm can adapt by making a different choice of the circulant.

We now prove that when A_{∞} is positive definite and *n* is sufficiently large, the circulants C_n are uniformly positive definite. Of course the finite sections A_n are also positive definite. The point is that an indefinite C_n —the possibility illustrated in Example 2—cannot continue as *n* increases.

THEOREM 1. Suppose $f(\theta) = \sum_{-\infty}^{\infty} a_k e^{ik\theta}$ is real and positive and in the Wiener class $(\sum |a_k| < \infty)$. Then the circulants C_n and C_n^{-1} are uniformly bounded and positive definite for large n.

Proof. The first column of C_n contains by construction the numbers $a_0, \dots, a_m, \dots, a_1$. (For simplicity we take *n* even and m = n/2.) Because the matrix is a circulant, its *j*th eigenvalue is

(3)
$$\lambda_i = a_0 + a_1 w^j + \dots + a_m w^{jm} + \dots + a_1 w^{j(n-1)}$$

Here $w = e^{2\pi i/n}$ is the primitive *n*th root of unity. The corresponding eigenvector of C (and of every circulant) is $x_j = (1, w^j, \dots, w^{j(n-1)})$; a direct multiplication gives $Cx_j = \lambda_j x_j$. Simplifying (3) by $w^n = 1$ yields

(3')
$$\lambda_{j} = a_{0} + a_{1}(w^{j} + w^{-j}) + \dots + a_{m-1}(w^{j(m-1)} + w^{j(1-m)}) + a_{m}w^{jm}.$$

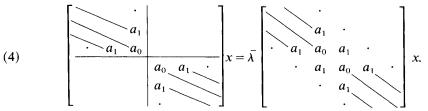
Thus the eigenvalue equals the partial sum from k = 1 - m to m of the series $\sum a_k e^{ik\theta}$, evaluated at the point $\theta = 2\pi j/n$ (where $e^{i\theta} = w^j$). Since the infinite series is absolutely convergent and its sum satisfies $f(\theta) \ge \delta > 0$, the partial sums are uniformly positive for large n and the proof is complete.

Example 3. Diagonally dominant matrices, with $a_0 > 2 \sum_{k \neq 0} |a_k|$, are positive definite and so are the circulants C_n .

Example 4. The function $f = \cosh \theta$ is even and positive. Therefore the matrices with $a_k = (1+k^2)^{-1} \cos k\pi$ lead to uniformly bounded C_n and C_n^{-1} , although the original matrix A is not diagonally dominant.

4. The limits of the eigenvalues. We come now to the central problem, to study the eigenvalues of $C_n^{-1}A_n$ for large *n*. In the next sections we transform that problem in order to carry out the analysis, and a Hankel matrix appears. At the end, when the limit is found, we transform back. The result was anticipated in [10], and it may be useful to separate its statement from the details of its proof.

THEOREM 2. As $n \to \infty$, the eigenvalues of $C_n^{-1}A_n$ approach the eigenvalues of the following doubly infinite problem:



The matrix on the right is an infinite circulant. The matrix on the left contains two back-to-back copies of the singly infinite Toeplitz matrix A_{∞} . Somehow the source of all the difficulty—the two boundaries that prevented the finite matrices A_n from being directly invertible by Fourier analysis—has reappeared in a new form.

Later in this paper we give an expression for the limiting eigenvalues $\overline{\lambda}$, by connecting them to a Hankel matrix and thus to a problem in rational approximation. That problem is approximation on the unit circle of a function $\tilde{v}(\theta)$ derived from $f(\theta)$, and it achieves our goal. The limits of the eigenvalues are determined from f. (The function has $|\tilde{v}(\theta)| \equiv 1$ and it appears as a "phase function" in systems theory [17]-[18] and apparently also in methods for numerical conformal mapping.)

At the end we return to the preconditioned conjugate gradient method, to prove superlinear convergence.

5. Orthogonal similarity and Hankel matrices. The key problem is $Ax = \lambda Cx$. There is a preliminary transformation which cuts this problem in half (for n = 2m), since all eigenvectors are odd or even:

(5)
$$x_{-} = \begin{bmatrix} y \\ -Jy \end{bmatrix}$$
 and $x_{+} = \begin{bmatrix} z \\ Jz \end{bmatrix}$ with $J = \begin{bmatrix} 1 \\ 1 \\ . \\ 1 \end{bmatrix}$

This property comes from the "centrosymmetry" of C and A, and it leads to the orthogonal transformation suggested by Cantoni and Butler [19]:

$$Q = \frac{1}{\sqrt{2}} \begin{bmatrix} I & I \\ -J & J \end{bmatrix}.$$

This combination of the identity I and the counteridentity J will produce two diagonal blocks in $Q^{-1}AQ$ and $Q^{-1}CQ$. Suppose we write

(6)
$$A = \begin{bmatrix} T & R \\ R^T & T \end{bmatrix} \text{ and } C = \begin{bmatrix} T & S \\ S^T & T \end{bmatrix},$$

in which T, R, S are Toeplitz of order m. T is symmetric around its main diagonal a_0 , and S is symmetric around a_m , while R has diagonals a_1, \dots, a_{n-1} and is displayed below. From JTJ = T and $JRJ = R^T$ we reach

(7)
$$Q^{-1}AQ = Q^{T}AQ = \begin{bmatrix} T - RJ & 0\\ 0 & T + RJ \end{bmatrix},$$
$$Q^{-1}CQ = Q^{T}CQ = \begin{bmatrix} T - SJ & 0\\ 0 & T + SJ \end{bmatrix}.$$

Thus the eigenvalue problem $Ax = \lambda Cx$ splits into

(8)
$$(T-RJ)y = \lambda_{-}(T-SJ)y \text{ and } (T+RJ)z = \lambda_{+}(T+SJ)z.$$

(*Note.* Those equations also appear directly when (5) and (6) are substituted into $Ax = \lambda Cx$.) There are *m* eigenvalues λ_+ and *m* eigenvalues λ_- which together represent the n = 2m eigenvalues λ . The eigenvectors x_+ and x_- are *C*-orthogonal as required $(x_+^T Cx_- = 0)$, and they are also orthogonal.

We emphasize the effect of the counteridentity J. The matrices RJ and SJ are no longer Toeplitz. Instead *they are Hankel matrices*. Like J itself, they are constant down each counterdiagonal. The *i*, *j* entry depends on the *sum* i+j instead of the difference i-j:

$$R = \begin{bmatrix} a_m & a_{n-1} \\ & \ddots & \\ a_1 & & a_m \end{bmatrix} \text{ and } RJ = \begin{bmatrix} a_{n-1} & a_m \\ & \ddots & \\ a_m & & a_1 \end{bmatrix}.$$

Thus equation (8) becomes a Toeplitz-Hankel eigenvalue problem. The northeast and southwest quarters of the original Toeplitz matrices have swung into blocks on the diagonal, and in the process they have become Hankel matrices.

A Hankel matrix is determined from its first column entries v_1, v_2, \dots by $V_{ij} = v_{i+j-1}$. In the singly infinite case its operator norm comes from the associated function $v(\theta) = \sum_{1}^{\infty} v_j e^{ij\theta}$, by Nehari's theorem [20]-[21], but not quite in the same way that the norm of a Toeplitz matrix comes from $f(\theta)$: The anti-analytic terms are at our disposal:

(9)
$$||A|| = \sup |f(\theta)|$$
 and $||V|| = \sup |\tilde{v}(\theta)| = \inf_{v_0, v_{-1}, \cdots, \theta} \left| \sum_{-\infty}^{\infty} v_j e^{ij\theta} \right|.$

We emphasize that the eigenvalues in the two cases are very different. The spectrum of A is the interval $[f_{\min}, f_{\max}]$, while V is a *compact operator* for v in l^1 —and its eigenvalues are the errors in a rational approximation problem.

6. The limiting equation $\overline{H}\overline{z} = \overline{\nu}\overline{T}\overline{z}$. This section begins our first approach to the eigenvalues of $C_n^{-1}A_n$ as $n \to \infty$. We return to (8), where the problem was split in half, and look at either of the two problems of order m:

(10)
$$(T+RJ)z = \lambda_+(T+SJ)z.$$

Our intention is to find the limits of the eigenvalues λ_+ as $n \to \infty$; numerical experiments indicate that limits exist.

First a simplification. Recall that the Hankel matrices RJ and SJ are identical on and below the main counterdiagonal. Their difference is the Hankel matrix:

$$H = SJ - RJ = \begin{bmatrix} h_1 & h_2 & \cdot & 0 \\ h_2 & & & 0 \\ \cdot & & & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \text{ with } h_j = a_j - a_{n-j} \text{ for } j \leq m.$$

Subtracting $\lambda_+(T+RJ)z$ from both sides of (10) and dividing through by λ_+ , we obtain the result

(11)
$$\frac{(1-\lambda_+)}{\lambda_+}(T+RJ)z = Hz.$$

We study that form of the problem, writing ν for the eigenvalue: $Hz = \nu (T + RJ)z$ with

$$\nu = \frac{1 - \lambda_+}{\lambda_+} \quad \text{and} \quad \lambda_+ = \frac{1}{1 + \nu}.$$

The clustering of λ_+ around 1 corresponds to the clustering of ν around 0.

We go directly to a statement of the limiting problem, and then consider its justification. As the order n increases, T and H approach singly infinite Toeplitz and Hankel matrices:

$$\bar{T} = \begin{bmatrix} a_0 & a_1 & a_2 & \cdot \\ a_1 & a_0 & a_1 & \cdot \\ a_2 & a_1 & a_0 & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix} \text{ and } \bar{H} = \begin{bmatrix} a_1 & a_2 & a_3 & \cdot \\ a_2 & a_3 & \cdot \\ a_3 & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$$

This is strong convergence of operators, if we think of all operators as acting on l^1 or l^2 . The matrices T and H represent operators T_m and H_m , which act as the zero operator after the first m components:

$$T_m = \begin{bmatrix} T & 0 \\ 0 & 0 \end{bmatrix} \text{ and } H_m = \begin{bmatrix} H & 0 \\ 0 & 0 \end{bmatrix}$$

LEMMA 1. The Hankel sequence H_m converges uniformly to $\overline{H} : \|\overline{H} - H_m\| \to 0$. The Toeplitz sequence T_m converges strongly to \overline{T} , and the sequence $(RJ)_m$ converges strongly to the zero operator:

$$\|\bar{T}x - T_m x\| \to 0$$
 and $\|\begin{bmatrix} RJ & 0\\ 0 & 0\end{bmatrix} x\| \to 0$ for each vector x .

Proof. For the Hankel matrices $\overline{H} = H_m$ the norm is given by (9):

(12)
$$\|\bar{H} - H_m\| \leq \sup \left| \sum_{1}^{m-1} a_{n-j} e^{ij\theta} + \sum_{m}^{\infty} a_j e^{ij\theta} \right|$$
$$\leq 2 \sum_{m}^{\infty} |a_j| \to 0 \quad \text{as } m \to \infty.$$

For the others, the estimate for $(RJ)_m$ is typical. With N fixed, the N-by-N submatrix in the upper left corner has *operator* norm going to zero. (Its entries are a_{n-j} , displayed earlier.) The larger entry a_1 is in the (m, m) position and it moves out as $m \to \infty$. There is convergence to zero for each fixed x but not uniform convergence—pointwise convergence but not norm convergence. Here we give only the simplest consequence for the eigenvalues.

THEOREM 3. Each eigenvalue $\bar{\nu}$ of the infinite Hankel-Toeplitz problem

$$\bar{H}\bar{z} = \bar{\nu}\bar{T}\bar{z}$$

is a limit of eigenvalues of the finite problems $Hz = \nu (T + RJ)z$.

Proof. The eigenvector \overline{z} will be our fixed vector x. Then the strong convergence noted above gives

$$\|H_m\bar{z}-\bar{\nu}(T+RJ)_m\bar{z}\|\to 0$$

Looking only at the *m*-by-*m* submatrices, where the operators are nonzero, and at the vector z_m taken from the first *m* components of \overline{z} , this is

(14)
$$||Hz_m - \bar{\nu}(T + RJ)z_m|| \to 0.$$

By writing B for T + RJ, this means that $||(H - \bar{\nu}B)^{-1}|| \to \infty$ as $m \to \infty$. In case $H - \bar{\nu}B$ is singular it means that $\bar{\nu}$ is an exact eigenvalue of the finite problem.

Because this is the generalized eigenvalue problem, with a matrix B on the right side instead of the identity, we need one extra step. Notice that these matrices B = T + RJ are *uniformly invertible*. They are diagonal blocks in the original Q^TAQ of (7). The matrices A are uniformly invertible because they are finite sections of a positive definite singly infinite Toeplitz matrix (it is \overline{T} !). Therefore we can convert to the ordinary eigenvalue problem for $B^{-1/2}HB^{-1/2}$, maintaining symmetry with the symmetric positive definite square root of B:

$$H - \bar{\nu}B = B^{1/2}(B^{-1/2}HB^{-1/2} - \bar{\nu}I)B^{1/2}.$$

In the l^2 matrix norm this yields

(15)
$$\|(H - \bar{\nu}B)^{-1}\| \le \|B^{-1}\| \max \frac{1}{|\nu_i - \bar{\nu}|}$$

Since $||B^{-1}||$ is bounded and the left side approaches infinity, we conclude that some eigenvalue ν_i of the finite problem converges to $\bar{\nu}$ as $m \to \infty$.

This completes an argument that could be made more precise. The fact that the eigenvectors belong to l^1 was established by Adamjan, Arov, and Krein [22] and pointed out to us by Nick Trefethen, who also led us into the eigenvalue theory of Hankel operators. In the Kac-Murdock-Szegö example $a_k = t^k$, where nearly half the eigenvalues are $\lambda = 1/(1+t^m)$, there is geometric convergence to $\overline{\lambda} = 1$. In terms of $\nu = (1-\lambda)/\lambda = t^m$, it is geometric convergence to the limit $\overline{\nu} = 0$. In that special example we will show that the only nonzero $\overline{\nu}$ is the number t, corresponding to the eigenvalue $\lambda = 1/(1+t)$ that stays away from 1.

We now check on the number of eigenvalues of the finite problem that are outside an interval around 1, after a remark on the other half of our original problem—the λ_{-} eigenvalues of $C^{-1}A$.

7. The twin problem $(T - RJ)y = \lambda_{-}(T - SJ)y$. The same simplification as in (11), but now subtracting $\lambda_{-}(T - RJ)y$ from both sides of the twin problem and dividing by λ_{-} , yields

(16)
$$\frac{(1-\lambda_{-})}{\lambda_{-}}(T-RJ)y = -Hy.$$

In this case we set

$$\mu = \frac{1 - \lambda_{-}}{\lambda_{-}}$$
 and, thus, $\lambda_{-} = \frac{1}{1 + \mu}$.

The finite problem is $Hy = -\mu (T - RJ)y$. Exactly as before, the strong convergence of H to \overline{H} , T to \overline{T} , and RJ to 0 leads to the limiting problem

$$\bar{H}\bar{y} = -\bar{\mu}\,\bar{T}\bar{y}.$$

This is identical to $\overline{H}\overline{z} = \overline{\nu}\overline{T}\overline{z}$ except for the sign change: $\overline{\mu} = -\overline{\nu}$. As before, each $\overline{\mu}$ is the limit of eigenvalues μ of the finite problem. There is an interesting corollary for the original eigenvalues λ_{-} and λ_{+} : In the limit there are pairs λ_{-} and λ_{+} which satisfy

(17)
$$\frac{1}{\lambda_{-}} + \frac{1}{\lambda_{+}} = 2$$

The left side approaches $(1 + \bar{\mu}) + (1 + \bar{\nu}) = 2$. This was first noticed by Alan Edelman in MATLAB experiments.

With this pair of limit problems, we have completed the proof of Theorem 2. The splitting into odd and even eigenvectors of that doubly infinite eigenvalue problem gives exactly (13) and its twin, with the same similarity Q and change from λ to ν and μ as in the finite case. A corresponding limit could be found for other constructions of the circulant C—and for multidimensional Toeplitz equations, in which our algorithm may be particularly useful. We concentrate here on understanding more clearly the asymptotic behavior for this choice of $C^{-1}A$.

8. The clustering of the spectrum of $C^{-1}A$. The limit problem gives us precise information about the asymptotic behavior of the algorithm. Even without that knowledge we can show that the eigenvalues of $C^{-1}A$ cluster at 1, by using the theory of collectively compact operators:

A family S of bounded operators on l^2 is collectively compact if the set $\{Kx: K \in S, \|x\| \le 1\}$ is relatively compact in l^2 .

This applies to $\{H_m\}$ and $\{\bar{H}-H_m\}$ when $\sum |h_j| < \infty$. For every ε we can choose $n(\varepsilon)$ so that outside the leading submatrix of that order, all these matrices have norms less than ε . Anselone [23] established the following consequences for approximation of the spectrum.

LEMMA 2 [23, Thm. 4.8]. If an open set contains the spectrum of \overline{H} , it also contains the spectrum of H_m for all sufficiently large m.

LEMMA 3 [23, Thm. 4.14]. The number of eigenvalues of H_m in a small ball around a nonzero eigenvalue μ of \overline{H} is, for large m, no greater than the multiplicity of μ .

Since \overline{H} is compact, those multiplicities are finite. (Our \overline{H} is the norm limit of the finite-dimensional H_m , by (12). From Hartman's theorem [21] it remains compact if its associated function is only continuous, and not necessarily in the Wiener class.) It follows easily that the eigenvalues of H_m are clustered around zero.

THEOREM 4. There exist $M(\varepsilon)$ and $N(\varepsilon)$ such that for m > M, at most N eigenvalues of H_m and of $(T \pm RJ)^{-1}H$ have absolute value exceeding ε .

Only a finite number of eigenvalues of \overline{H} have $|\mu| > \varepsilon$. Therefore Lemma 2 allows us to locate the eigenvalues of H_m for large m, and Lemma 3 allows us to count them. The total cannot exceed the total for \overline{H} .

It is this count that was not included in Theorem 3 on the limiting values. As there, we have to handle the matrices $B = T \pm RJ$ on the right side of the eigenvalue problem. They are diagonal blocks in $Q^T A Q$, and their eigenvalues are between f_{min} and f_{max} (where $f = \sum a_k e^{ik\theta}$). Therefore the eigenvalues of $(T \pm RJ)^{-1}H$ are also counted by Theorem 4.

Thus the eigenvalues λ of $C^{-1}A$ cluster around 1. Only a fixed number, independent of *n*, can be outside $(1 - \varepsilon, 1 + \varepsilon)$. Now we go back to the equation $\overline{H}\overline{z} = \overline{\nu}\overline{T}\overline{z}$, which identifies the asymptotic limits of the eigenvalues.

9. The eigenvalues of $\overline{T}^{-1}\overline{H}$. It is this singly infinite problem that is attractive to work with, because all the information is in the function $f(\theta) = \sum a_k e^{ik\theta}$. The difficulty is to extract it. One preliminary difficulty is that we still have a generalized eigenvalue problem, with \overline{T} on the right-hand side: $\overline{H}\overline{z} = \overline{\nu}\overline{T}\overline{z}$. The inverse of \overline{T} is not Toeplitz, and $\overline{T}^{-1}\overline{H}$ is not Hankel, but there is a way to preserve those properties, by factoring \overline{T} and putting part of \overline{T}^{-1} on each side of \overline{H} :

 $\overline{T} = WW^T$ with W = upper triangular Toeplitz operator.

This is equivalent to representing the positive function f as a square:

$$f = \sum_{-\infty}^{\infty} a_k e^{ik\theta} = \left| \sum_{-\infty}^{0} w_k e^{ik\theta} \right|^2 = |w|^2.$$

From the Wiener theory the sequence \cdots , w_{-2} , w_{-1} , w_0 is in l^1 . The function w, with no zeros on or outside the unit circle, corresponds to W in the same way that f corresponds to \overline{T} :

$$w = \sum_{-\infty}^{0} w_k e^{ik\theta} \leftrightarrow W = \begin{bmatrix} w_0 & w_{-1} & w_{-2} & \cdot \\ & w_0 & w_{-1} & \cdot \\ & & w_0 & \cdot \\ & & & & \cdot \end{bmatrix}.$$

Note that w is *anti-analytic*, with negative k. The same properties hold for the matrix $U = W^{-1}$, associated with the reciprocal function $u = w^{-1}$:

$$u(\theta) = \sum_{-\infty}^{0} u_k e^{ik\theta} = \left(\sum_{-\infty}^{0} w_k e^{ik\theta}\right)^{-1}.$$

These functions will take us from the Toeplitz-Hankel product $\overline{T}^{-1}\overline{H}$ (the eigenvalue problem with \overline{T} on the right side) to a single Hankel matrix V. It has the same eigenvalues $\overline{\nu}$, and to study them we need to know its associated function.

THEOREM 5. The matrix $\overline{T}^{-1}\overline{H} = W^{-T}W^{-1}\overline{H} = U^{T}UH$ is similar to the Hankel matrix $V = UHU^{T}$. The associated function is

(18)
$$v(z) = \sum_{1}^{\infty} v_k z^k = analytic \ part \ of \ \frac{\sum_{1}^{\infty} a_k z^k}{(\sum_{-\infty}^{\infty} w_k z^k)^2}$$
$$= analytic \ part \ of \ \bar{w}/w.$$

Proof. Certainly $U^T UH$ is similar to $V = UHU^T$. To verify that this is a Hankel matrix, and to identify its function v(z), we carry out an example with two nonzero coefficients h_1 , h_2 and u_0 , u_1 . (The general rule is that upper triangular Toeplitz times Hankel is Hankel, and Hankel times lower triangular Toeplitz is Hankel; we hope the example will be convincing.) The infinite matrices can be cut off after three rows and columns:

$$\begin{bmatrix} u_0 & u_1 & 0 \\ 0 & u_0 & u_1 \\ 0 & 0 & u_0 \end{bmatrix} \begin{bmatrix} h_1 & h_2 & 0 \\ h_2 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} u_0 & 0 & 0 \\ u_1 & u_0 & 0 \\ 0 & u_1 & u_0 \end{bmatrix} = \begin{bmatrix} h_1 u_0^2 + 2h_2 u_1 u_0 & h_2 u_0^2 & 0 \\ h_2 u_0^2 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

The corresponding multiplication of functions correctly gives the analytic part of $(h_1z + h_2z^2)(u_0 + u_1/z)^2 = h_1u_0^2z + 2h_2u_0u_1z + h_2u_0^2z^2 = v(z)$. This verifies the first line of (18), since u = 1/w, and we only mention a surprise in comparison with the Toeplitz case. If the matrix H in the middle were Toeplitz, with associated function h, then the product UHU^T would again be Toeplitz, but its associated function would be $h|u|^2$ —where in the Hankel case it is the analytic part of hu^2 . Note that the analytic part is taken to start at the linear term in z— because the Hankel matrix starts at v_1 .

Now we look at the last step in (18), the elegant form for v. It was noticed by Alan Edelman, and the second author observed that it follows immediately from the line above: We can add the anti-analytic terms $\sum_{-\infty}^{0} a_k z^k$ to the numerator to obtain f, without changing the analytic part. Therefore

v is the analytic part of
$$\frac{f}{w^2} = \frac{w\bar{w}}{w^2} = \frac{\bar{w}}{w}$$
.

Now the asymptotic problem is the spectrum of V. We recall that it is a compact operator, so its eigenvalues $\bar{\nu}$ cluster at zero. It is also a Hankel operator, and the eigenvalues are connected to a part of mathematics that looks entirely separate: approximation by rational functions on the unit circle.

10. Hankel eigenvalues and rational approximation. We recall the main facts from [22] and [24]. The singular values of a Hankel matrix V—which in our symmetric case means the absolute values $|\bar{\nu}_0| \ge |\bar{\nu}_1| \ge \cdots$ —are the errors in the best approximation of the function v(z) from the class \tilde{R}_n . The error is measured in the sup norm on the unit circle:

(19)
$$\left\|\bar{\nu}_{n}\right\| = \inf_{r \in \tilde{R}_{+}} \|v(z) - \tilde{r}(z)\|_{\infty}.$$

It is important that \tilde{R}_n is larger than the class R_n of rational functions (ratios of polynomials of degree *n*). The anti-analytic part of \tilde{r} is arbitrary: $\tilde{R}_n = R_n + anti-analytic$.

Thus the approximation problem for v is the same as for \bar{w}/w , because the analytic parts are the same. The approximants have the form

$$\tilde{r} = \sum_{-\infty}^{n} d_{k} z^{k} / \sum_{0}^{n} e_{k} z^{k},$$

where the numerator is bounded in any bounded subset of $\{|z| \ge 1\}$. The optimal error curve $v - \tilde{r}^*$ is a circle of radius $|\bar{v}_n|$ around the origin. Except in degenerate cases its winding number is 2n+1 [25].

Thus the estimation of the asymptotic eigenvalues $\bar{\nu}$ is equivalent to a problem in rational approximation. That may not be a simplification; most applications go the other way. Of course there is a special (but important) case when f itself is rational; this corresponds to banded Toeplitz matrices times banded inverses, and the Hankel matrix V has only finitely many nonzero eigenvalues. That was the case for the Kac-Murdock-Szegö example in which

(20)
$$f = \frac{1 - t^2}{|1 - tz|^2}, \quad u = \frac{1 - t/z}{(1 - t^2)^{1/2}}, \quad v = \frac{(1 - t^2)tz}{1 - tz}$$

Here v maps the unit circle to a circle with center at t^2 and radius t. Therefore the best constant approximation to v is $r = t^2$, and the error is $|\bar{v}| = t$. The corresponding λ 's are $1/(1 \pm t)$. That is the correct limit, from our explicit calculation, of the original eigenvalue problem $Ax = \lambda Cx$, in which all other eigenvalues approached 1.

A second example will bring out the important point, which is the very rapid decrease of the errors $|\bar{v}|$ in rational approximation. Suppose the matrix is not tridiagonal but pentadiagonal. The function f will be

$$f = |(1 - t e^{i\theta})(1 - s e^{i\theta})|^2.$$

As s approaches zero this goes back to the previous example (or more precisely to its inverse—it is a corollary of (18) that f and f^{-1} lead to the same results, and now it is A instead of A^{-1} that is banded). There are *two* nonzero errors $|\bar{\nu}_0|$ and $|\bar{\nu}_1|$, after which the approximation is exact and all eigenvalues approach $\bar{\nu} = 0$ (which is $\lambda = 1$). Those two errors are given by the quadratic equation

(21)
$$\nu^2 - \nu(t+s)(ts-1) - t^2 s^2 = 0.$$

For $s = t = \frac{1}{2}$ the limiting eigenvalues are $|\bar{\nu}_0| = .825$ and $|\bar{\nu}_1| = .076$. Thus it is not the case that the two limits $|\bar{\nu}|$ are near t and s. For t = s the leading terms in $|\bar{\nu}|$ are 2t and $t^3/2$, and it is this *cube* of t that indicates rapid decrease. A similar phenomenon was noticed by Trefethen [24, Thm. 6.3].

We turn to the consequences for the iterative algorithm when there is rapid decrease of the $|\bar{\nu}|$.

11. The rate of convergence to $x = A^{-1}b$. The conjugate gradient method is a recursive calculation of a sequence of projections. After q cycles, x_q is as close as possible (in an appropriate norm) to the solution $x = A^{-1}b$, among all vectors in the Krylov subspace spanned by $C^{-1}b$, $C^{-1}AC^{-1}b$, $C^{-1}AC^{-1}b$, \cdots . This makes possible an estimate of the error $e_q = x - x_q$:

(22)
$$\|e_q\| \leq [\min_{P_q} \max_{\lambda} |P_q(\lambda)|] \|e_0\|.$$

The maximum is taken over the eigenvalues of $C^{-1}A$. The minimum is over polynomials of degree q with constant term 1. The problem is to estimate that minimum.

One choice of P_q , if the eigenvalues are known to lie in the interval $[\alpha, \beta]$, is the Chebyshev choice: the best polynomial when the maximization is taken over all $\alpha \leq \lambda \leq \beta$. That has the drawback of using only α and β ; it cannot take advantage of clustering of the eigenvalues. By scale invariance, the estimate depends only on the condition number β/α .

At the other extreme, we can choose P_q to annihilate the q extreme eigenvalues. In our problem those eigenvalues come in pairs λ_+ and λ_- , on opposite sides of 1, and such a pair is annihilated by the factor

(23)
$$p(x) = \left(1 - \frac{x}{\lambda_+}\right) \left(1 - \frac{x}{\lambda_-}\right).$$

Between those roots the maximum of |p| is attained at the average $x = \frac{1}{2}(\lambda_+ + \lambda_-)$, where $|p| = (\lambda_+ - \lambda_-)^2/4\lambda_+\lambda_-$. It is easy to find the asymptotic convergence rate of the conjugate gradient method in the important case when the rational approximation errors decrease geometrically to zero:

(24)
$$|\bar{\nu}_i| = O(r^j)$$
 with $r < 1$.

THEOREM 6. Suppose (24) holds, which depends on the original f. (It is certainly true if f is analytic in a neighborhood of |z|=1.) Then the errors in the circulant-preconditioned conjugate gradient method decrease, asymptotically as $n \to \infty$, at the superlinear rate

(25)
$$||e_q|| \leq c^q r^{q^2/4+q/2} ||e_0||.$$

The decisive factor is $r^{q^2/4}$. To find it we note that $\lambda_{\pm} = 1/(1 \mp \bar{\nu})$:

$$|p_j| \leq \frac{(\lambda_{+j} - \lambda_{-j})^2}{4\lambda_{+j}\lambda_{-j}} = \frac{\overline{\nu}_j^2}{1 - \overline{\nu}_j^2} \leq c^2 r^{2j} \quad \text{by (24).}$$

The polynomial $P_{2q} = p_1 p_2 \cdots p_q$ annihilates the q extreme pairs of eigenvalues and satisfies

$$|P_{2q}(\lambda)| \leq c^{2q} r^2 r^4 \cdots r^{2q} = c^{2q} r^{q^2+q}.$$

This holds for all λ in the inner interval between λ_{+q} and λ_{-q} , where the remaining eigenvalues are. Therefore (25) comes directly from (22), after a change from 2q to q.

12. Superlinear convergence. In a sense the convergence rate $r^{q^2/4}$ was the object of this paper. By connecting the eigenvalues of $C^{-1}A$ to the spectrum of the Hankel matrix V, and by assuming (24), we were led to that unusual rate. Trefethen observed that (24) will hold *if the original* $f(z) = \sum a_j z^j$ *is analytic in a neighborhood of* |z| = 1, and, further, that this condition is *far from necessary*. We expect that rate for a wide class of applications, but we can also prove superlinear convergence in its absence.

For this we modify the polynomial P_q as in [26]-[27]. It will annihilate N pairs of extreme eigenvalues, by including N of the quadratic factors (23). The remaining factor of degree q-2N will be the Chebyshev choice on the interval between λ_{-N} and λ_{+N} , which contains the remaining eigenvalues. We know from Theorem 4 (which applied to all f in the Wiener class) that for any fixed ε , all but $N(\varepsilon)$ eigenvalues are within ε of 1. The N extreme eigenvalues are in the fixed interval

$$\frac{f_{\min}}{f_{\max}} - \varepsilon < \lambda (C^{-1}A) < \frac{f_{\max}}{f_{\min}} + \varepsilon$$

for large *n*, applying Theorem 1 to C^{-1} and the elementary bounds $f_{\min} \leq \lambda(A) \leq f_{\max}$ to *A*. The quadratics *p* that annihilate those extreme pairs are bounded by a fixed constant *K* on $(1 - \varepsilon, 1 + \varepsilon)$. The other (Chebyshev) factor of degree q - 2N is of order ε^{q-2N} on that interval. Therefore the polynomial P_q satisfies a crude bound

$$|P_{q}(\lambda)| \leq c \varepsilon^{q-2N} K^{N} \leq C(\varepsilon) \varepsilon^{q}$$

for all eigenvalues λ of $C^{-1}A$, when the order *n* is sufficiently large.

It follows from (22) that $||e_q|| \leq C(\varepsilon)\varepsilon^q ||e_0||$. The number of iterations to achieve a fixed accuracy remains bounded as the matrix order *n* is increased. Each iteration requires $O(n \log n)$ operations using the FFT. Therefore the work to obtain the solution $x = A^{-1}b$ to given accuracy δ is $c(f, \delta)n \log n$. The real question is the efficiency in practice, and we will be glad to see this straightforward algorithm tried on genuine applications.

13. Tentative experiments. One family of Toeplitz matrices is particularly convenient for testing, and we report here on the results. The entries down the kth diagonal (with k = 1 for the main diagonal) are $a_k = k^{-p}$. For p = 2 the entries $1, \frac{1}{4}, \frac{1}{9}, \cdots$ decrease quickly away from the center. At p = 1 the sum $1 + \frac{1}{2} + \frac{1}{3} + \cdots$ is divergent, and we leave the Wiener class. At $p = \frac{1}{2}$ the eigenvalues are less clustered (and we had not known that A and C were positive definite). But even at p = .01 the conjugate gradient convergence is remarkable. Nine or ten steps reduce the residual by 10^{-8} , independent of the order n. That last statement is entirely experimental (this is MATLAB mathematics), because the theory for the Wiener class has been left behind.

Table 1 shows the four largest eigenvalues of $C^{-1}A$ when n = 40. We see how the largest grows as p decreases (and the smallest is related to it by $\lambda_{\min}^{-1} + \lambda_{\max}^{-1} \approx 2$). However, there is still strong clustering around $\lambda = 1$.

Table 2 shows the norms of the residuals $r_q = b - Ax_q$ in the conjugate gradient method. It is preconditioned by the circulant *C*, and the right-hand side *b* has randomly chosen entries from a uniform distribution over (0, 1). *N* is the number of iterations to reach a residual norm below 10^{-8} . In these four cases, the smallest eigenvalues of *C* were, respectively, .645, .385, .207, and .004. Alan Edelman convinced us that the asymptotic behavior of this eigenvalue (small *p* and large *n*) is $(\log \pi/2)p$. The positive definiteness is significant; we give an indefinite example next.

	IABLE I							
р	λ1	λ ₂	λ ₃	λ_4				
2	1.360	1.029	1.003	1.002				
1	2.072	1.079	1.018	1.013				
$\frac{1}{2}$	3.100	1.111	1.049	1.035				
$\frac{1}{100}$	5.596	1.190	1.136	1.102				

TABLE 2

p	$\ r_1\ $	$\ r_2\ $	$\ r_3\ $	$\ r_4\ $	N
2	3.441	.094	.007	.000	6
1	4.131	.338	.031	.014	7
$\frac{1}{2}$	4.136	.634	.290	.028	8
$\frac{1}{100}$	3.898	2.560	.219	.027	10

Suppose that the Toeplitz entries are $a_k = 1/k!$. The underlying function is

$$f(z) = e^{z} + e^{1/z} - 1$$

and there will be extremely good approximation by rational functions, which implies quick convergence of the Hankel eigenvalues $\nu \to 0$, and strong clustering of the Toeplitz-circulant eigenvalues $\lambda \to 1$. Analyticity is clear, but positive definiteness is not. The function f is negative at z = -1, and $\lambda_{\min}(C) = -.264$ with n = 40. Nevertheless, the eigenvalues of $C^{-1}A$ cluster near 1 (amazingly so, since some are complex). What may be of interest is the sequence of norms of residuals, which show fast convergence after a kink to digest the indefiniteness of the problem:

$$\|r\| = 3.59$$
 .51 1.38 2.29 .22 .01 .00.

By adding 1 to the main diagonal and changing f to $e^{z} + e^{1/z}$, the matrices A and C become positive definite (though still not diagonally dominant). $C^{-1}A$ has only three eigenvalues that are further than 10^{-5} above $\lambda = 1$, and three more at that distance below:

$$\lambda(C^{-1}A) = 2.02$$
 1.06 1.0009 1.000007 1.00000003.

The result is convergence in six steps of the preconditioned conjugate gradient method.

In all these cases the operation count is $O(n \log n)$.

REFERENCES

- N. LEVINSON, The Wiener RMS (Root-Mean Square) error criterion in filter design and prediction, J. Math. Phys., 25 (1947), pp. 261-278.
- [2] G. AMMAR AND W. GRAGG, Superfast solution of real positive definite Toeplitz systems, SIAM J. Matrix Anal. Appl., 9 (1988), pp. 61–76.
- [3] F. DE HOOG, A new algorithm for solving Toeplitz systems of equations, Linear Algebra Appl., 88/89 (1987), pp. 122-138.
- [4] J. BUNCH, Stability of methods for solving Toeplitz systems of equations, SIAM J. Sci. Statist. Comput., 6 (1985), pp. 349-364.
- [5] T. KAILATH, A theorem of I. Schur and its impact on modern signal processing, in I. Schur: Methods in Operator Theory and Signal Processing, I. Gohberg, ed., Birkhäuser, Basel, 1986.
- [6] R. BRENT AND F. LUK, A systolic array for the linear-time solution of Toeplitz systems of equations, J. VLSI Comput. Syst., 1 (1983), pp. 1–22.
- [7] G. STRANG, A proposal for Toeplitz matrix calculations, Stud. Appl. Math., 74 (1986), pp. 171-176.
- [8] P. DELSARTE AND Y. GENIN, The split Levinson algorithm, IEEE Trans. Acoust. Speech Signal Process., 34 (1986), pp. 470-478.
- [9] G. STRANG, Introduction to Applied Mathematics, Wellesley-Cambridge Press, Wellesley, MA, 1986.
- [10] G. STRANG AND A. EDELMAN, The Toeplitz-circulant eigenvalue problem $Ax = \lambda Cx$, Oakland Conference on PDE's, L. Bragg and J. Dettman, eds., Longmans, London, 1987.
- [11] M. KAC, W. MURDOCK, AND G. SZEGÖ, On the eigenvalues of certain Hermitian forms, Arch. Rational Mech. Anal., 2 (1953), pp. 767-800.
- [12] P. DELSARTE AND Y. GENIN, Spectral properties of finite Toeplitz matrices, Philips Research Lab. Manuscript M50, Proc. Internat. Symposium Math. Theory of Networks and Systems, Beer-Sheva, Israel, 1983.
- [13] G. CYBENKO, On the eigenstructure of Toeplitz matrices, IEEE Trans. Acoust. Speech Signal Process., 32 (1984), pp. 918-920.
- [14] W. F. TRENCH, On the eigenvalue problem for Toeplitz matrices generated by rational functions, Linear and Multilinear Algebra, 17 (1985), pp. 337-353.
- [15] U. GRENANDER AND G. SZEGÖ, Toeplitz Forms and Their Applications, University of California Press, Berkeley, CA, 1958.
- [16] R. M. GRAY, On the asymptotic eigenvalue distribution of Toeplitz matrices, IEEE Trans. Inform. Theory, 18 (1972), pp. 725-729.
- [17] M. GREEN AND B. D. O. ANDERSON, The approximation of power spectra by phase matching, Proc. IEEE Conference on Decision and Control, Athens, GA, 1986.

- [18] K. GLOVER, R. CURTAIN, AND J. PARTINGTON, Realisation and approximation of linear infinite dimensional systems with error bounds, # 258, Engineering Department, Cambridge University, Cambridge, UK.
- [19] A. CANTONI AND P. BUTLER, Eigenvalues and eigenvectors of symmetric centrosymmetric matrices, Linear Algebra Appl., 13 (1976), pp. 275–288.
- [20] Z. NEHARI, On bounded bilinear forms, Ann. Math., 65 (1957), pp. 153-162.
- [21] S. C. POWER, Hankel Operators on Hilbert Space, Pitman, Boston, 1982.
- [22a] V. ADAMJAN, D. AROV, AND M. KREIN, Analytic properties of Schmidt pairs for a Hankel operator and the generalized Schur-Takagi problem, Math. USSR-Sb., 15 (1971), pp. 31-73.
- [22b] —, Infinite Hankel matrices and generalized Carathéodory-Féjér and Riesz problems, Functional Anal. Appl., 2 (1968), pp. 1–18.
- [23] P. M. ANSELONE, Collectively Compact Operator Approximation Theory, Prentice-Hall, Englewood Cliffs, NJ, 1971.
- [24] L. N. TREFETHEN, Rational Chebyshev approximation on the unit disk, Numer. Math., 37 (1981), pp. 297-320.
- [25] E. HAYASHI, L. N. TREFETHEN, AND M. GUTKNECHT, *The* CF *table*, Constructive Approximation, submitted.
- [26] H. VAN DER VORST AND A. VAN DER SLUIS, *The rate of convergence of conjugate gradients*, Preprint 354, University of Utrecht, Utrecht, the Netherlands, 1984.
- [27] O. AXELSSON AND G. LINDSKOG, On the rate of convergence of the preconditioned conjugate gradient method, Numer. Math., 48 (1986), pp. 499-523.