

# Optimal trigonometric preconditioners for nonsymmetric Toeplitz systems

Daniel Potts and Gabriele Steidl

220/96

**Abstract.** This paper is concerned with the solution of systems of linear equations  $\mathbf{T}_N \mathbf{x}_N = \mathbf{b}_N$ , where  $\{\mathbf{T}_N\}_{N \in \mathbb{N}}$  denotes a sequence of nonsingular nonsymmetric Toeplitz matrices arising from a generating function of the Wiener class. We present a technique for the fast construction of optimal trigonometric preconditioners  $\mathbf{M}_N = \mathbf{M}_N(\mathbf{T}'_N \mathbf{T}_N)$  of the corresponding normal equation. Moreover, we prove that the spectrum of the preconditioned matrix  $\mathbf{M}_N^{-1} \mathbf{T}'_N \mathbf{T}_N$  is clustered at 1 such that the CG-method applied to the normal equation converges superlinearly. Numerical tests confirm the theoretical expectations.

1991 *Mathematics Subject Classification.* 65F10, 65F15, 65T10.

*Key words and phrases.* Toeplitz matrix, Krylov space methods, CG-method, preconditioners, normal equation, clusters of eigenvalues.

# 1 Introduction

Consider the system of linear equations

$$\mathbf{T}_N \mathbf{x}_N = \mathbf{b}_N, \quad (1.1)$$

where  $\mathbf{T}_N \in \mathbb{R}^{N,N}$  denotes a nonsingular Toeplitz matrix. Toeplitz systems arise in a variety of applications in mathematics and engineering (see [7] and the references therein). While there exist fast direct Toeplitz solvers for Hermitian positive definite Toeplitz matrices  $\mathbf{T}_N$ , such techniques are not available in the non-Hermitian case. Iterative methods like GMRES and CG often provide a fast solution of (1.1) if they are applied in connection with preconditioning techniques [7]. In particular, these methods profit from the fact that the vector multiplication with the Toeplitz matrix  $\mathbf{T}_N$  in each iteration step can be computed with  $O(N \log N)$  arithmetical operations by using the fast Fourier transform (FFT). Clearly, the multiplication with the preconditioned matrix should have the same arithmetic complexity. Two types of preconditioners are mainly exploited for linear Toeplitz systems, namely optimal (Cesáro) circulant preconditioners  $\mathbf{M}_N = \mathbf{C}_N(\mathbf{T}_N)$  [5] and more simple so-called "Strang" circulant preconditioners  $\mathbf{M}_N = \mathbf{S}_N(\mathbf{T}_N)$  [6]. One reason for the choice of circulant preconditioners is the fact that circulant matrices can be diagonalized by the Fourier matrix  $\mathbf{F}_N := (e^{-2\pi i j k / N})_{j,k=0}^{N-1}$ , where the multiplication of a vector with  $\mathbf{F}_N$  takes only  $O(N \log N)$  arithmetical operations. Moreover, under certain assumptions on the generating function of  $\mathbf{T}_N$  (see [8], [22]), it can be proved that the singular values of  $\mathbf{M}_N^{-1} \mathbf{T}_N$  are clustered at 1. For non-Hermitian  $\mathbf{T}_N$ , this results in a superlinear convergence of the CG-method applied to the system

$$(\mathbf{M}_N^{-1} \mathbf{T}_N)^* (\mathbf{M}_N^{-1} \mathbf{T}_N) \mathbf{x}_N = (\mathbf{M}_N^{-1} \mathbf{T}_N)^* \mathbf{M}_N^{-1} \mathbf{b}_N. \quad (1.2)$$

To our knowledge, up to now, for non-Hermitian Toeplitz systems and Toeplitz least square problems, only circulant preconditioners with respect to Toeplitz matrices were constructed and used in some kind of normal equation as in (1.2) [8], [11] or as so-called displacement preconditioners [12]. When we finished the paper, we became aware of new results of E.E. Tyrtyshnikov et al. concerning the convergence behaviour of the preconditioned GMRES-method [23] which avoids the transition of (1.1) to the normal equation. However, the preconditioners are again (improved) circulants, which were constructed with respect to  $\mathbf{T}_N$ .

In this paper, we restrict our attention to nonsymmetric *real* Toeplitz matrices  $\mathbf{T}_N$ . Here, it seems to be natural, to replace the circulant matrices by matrices which are diagonalizable by some real trigonometric matrices. Of course, the commonly used trigonometric transforms are closely related to the Fourier transform. Indeed, for symmetric Toeplitz matrices  $\mathbf{T}_N$  with positive continuous  $2\pi$ -periodic generating functions, trigonometric preconditioning significantly accelerates the convergence of the CG-method.

In this paper, we suggest the solution of (1.1) by applying the CG-method to the preconditioned normal equation

$$\mathbf{M}_N^{-1} \mathbf{T}'_N \mathbf{T}_N \mathbf{x}_N = \mathbf{M}_N^{-1} \mathbf{T}'_N \mathbf{b}_N,$$

2

Zugangsnummer: 54194

Signatur:

UNIVERSITÄT MANNHEIM  
Fachsbibliothek Mathematik und Informatik

where in contrast to (1.2),  $\mathbf{M}_N = \mathbf{M}_N(\mathbf{T}'_N \mathbf{T}_N)$  denotes the optimal preconditioner with respect to  $\mathbf{T}'_N \mathbf{T}_N$ . We demonstrate that the construction of such optimal preconditioners can be realized with only  $O(N \log N)$  arithmetical operations despite the fact that  $\mathbf{T}'_N \mathbf{T}_N$  is no longer a Toeplitz matrix. We prove that under certain assumptions on  $\mathbf{T}_N$  the eigenvalues of  $\mathbf{M}_N^{-1} \mathbf{T}'_N \mathbf{T}_N$  are clustered at 1. Although our approach works in exactly the same way for different trigonometric transforms, we prefer to investigate the DCT-II preconditioner in detail and add only few facts concerning the other trigonometric preconditioners. We hope that our notation makes the approach for other trigonometric preconditioners immediately clear. Numerical tests were performed for the different trigonometric preconditioners. Note that in all examples, our preconditioning was superior over the method (1.2) with an optimal trigonometric preconditioner  $\mathbf{M}_N(\mathbf{T}_N)$  of  $\mathbf{T}_N$ .

This paper is organized as follows: Section 2 contains the basic matrix notation. In Section 3, we study the relations between trigonometric transforms and Toeplitz matrices. In particular, we introduce a method for the fast vector multiplication with real nonsymmetric Toeplitz matrices based on real trigonometric transforms. In Section 4, we introduce optimal trigonometric preconditioners. Section 5 is concerned with the proof that the eigenvalues of the preconditioned matrix  $\mathbf{M}_N^{-1} \mathbf{T}'_N \mathbf{T}_N$  are clustered at 1. In Section 6, we present the fast construction of the optimal preconditioner. Finally, Section 7 confirms the theoretical expectations by numerical tests.

## 2 Notation

For the sake of clarity, we collect the matrix notation in this preliminary section.

Let  $\mathbf{a}_N := (a_0, \dots, a_{N-1})'$ ,  $\mathbf{b}_N := (b_0, \dots, b_{N-1})'$  and let  $\mathbf{0}_N$  be the vector consisting of  $N$  zeros. Here  $\mathbf{A}'$  is the transpose of  $\mathbf{A}$ . By  $\mathbf{I}_N$  we denote the  $(N, N)$ -identity matrix and by  $\mathbf{e}_k \in \mathbb{R}^N$  the  $k$ -th identity vektor. To describe Toeplitz and Hankel matrices, we use the following notation:

$$\text{toeplitz}(\mathbf{a}', \mathbf{b}') := \begin{pmatrix} a_0 & a_1 & \dots & a_{N-1} \\ b_1 & a_0 & \dots & a_{N-2} \\ \vdots & \vdots & \ddots & \vdots \\ b_{N-2} & b_{N-3} & \dots & a_1 \\ b_{N-1} & b_{N-2} & \dots & a_0 \end{pmatrix} \quad (\text{with } a_0 = b_0),$$

stoepiltz  $\mathbf{a}'$ : symmetric Toeplitz matrix with first row  $\mathbf{a}'$ ,

atoeplitz  $\mathbf{a}'$ : antisymmetric Toeplitz matrix with first row  $\mathbf{a}'$ , where  $a_0 = 0$ ,

$$\text{hankel}(\mathbf{a}', \mathbf{b}') := \begin{pmatrix} a_0 & \dots & a_{N-2} & a_{N-1} \\ a_1 & \dots & a_{N-1} & b_{N-2} \\ \vdots & \vdots & \vdots & \vdots \\ a_{N-2} & \dots & b_2 & b_1 \\ a_{N-1} & \dots & b_1 & b_0 \end{pmatrix} \quad (\text{with } a_{N-1} = b_{N-1}),$$

shankel  $\mathbf{a}'$ : persymmetric Hankel matrix with first row  $\mathbf{a}'$ ,

ahankel  $\mathbf{a}'$ : antipersymmetric Hankel matrix with first row  $\mathbf{a}'$ , where  $a_{N-1} = 0$ .

Further, we introduce the matrices

$$\mathbf{Z}'_{N,1} := \begin{pmatrix} 0 & 1 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix} \in \mathbb{R}^{N,N+1}, \quad \mathbf{Z}'_{N,2} = \begin{pmatrix} 1 & \dots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & 1 & 0 \end{pmatrix} \in \mathbb{R}^{N,N+1},$$

and

$$\mathbf{R}'_N := \begin{pmatrix} 0 & 1 & & 0 & 0 \\ \vdots & & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & 0 \end{pmatrix} \in \mathbb{R}^{N-1,N+1}.$$

Let  $\text{diag } \mathbf{a}$  be the diagonal matrix with diagonal  $\mathbf{a}$  and let  $\delta(\mathbf{A}) := \text{diag}(a_{k,k})_{k=0}^{N-1}$ , where  $a_{k,k}$  is the  $(k, k)$ -th entry of  $\mathbf{A}$ . By

$$\text{tr } \mathbf{A} := \sum_{k=0}^{N-1} a_{k,k}$$

we denote the *trace* of  $\mathbf{A}$ . Moreover, we need the following *matrix norms*:

*Spectral norm*:

$$\|\mathbf{A}\|_2 := (\text{maximum of the absolute values of the singular values of } \mathbf{A})^{1/2},$$

*Frobenius norm*:

$$\|\mathbf{A}\|_F := \left( \sum_{j,k=0}^{N-1} a_{jk}^2 \right)^{1/2},$$

*1-norm*:

$$\|\mathbf{A}\|_1 := \max \left\{ \sum_{j=0}^{N-1} a_{jk} : k = 0, \dots, N-1 \right\}.$$

If it does not make confusion, we use the same notation for the norm of absolute summable sequences  $a = \{a_k\}_{k \in \mathbb{Z}} \in l_1$ , i.e.

$$\|a\|_1 := \sum_{k \in \mathbb{Z}} |a_k|.$$

### 3 Trigonometric transforms and Toeplitz matrices

We introduce four discrete sine transforms (DST) and four discrete cosine transforms (DCT) as classified by Wang [24]:

$$\text{DCT-I} : \quad \mathbf{C}_{N+1}^I := \left( \frac{2}{N} \right)^{1/2} \left( \varepsilon_j^N \varepsilon_k^N \cos \frac{jk\pi}{N} \right)_{j,k=0}^N \in \mathbb{R}^{N+1,N+1},$$

$$\text{DCT-II} : \quad \mathbf{C}_N^{II} := \left( \frac{2}{N} \right)^{1/2} \left( \varepsilon_j^N \cos \frac{j(2k+1)\pi}{2N} \right)_{j,k=0}^{N-1} \in \mathbb{R}^{N,N},$$

$$\text{DCT-III} : \quad \mathbf{C}_N^{III} := (\mathbf{C}_N^{II})' \in \mathbb{R}^{N,N},$$

$$\text{DCT-IV} : \quad \mathbf{C}_N^{IV} := \left( \frac{2}{N} \right)^{1/2} \left( \cos \frac{(2j+1)(2k+1)\pi}{4N} \right)_{j,k=0}^{N-1} \in \mathbb{R}^{N,N},$$

and

$$\begin{aligned}
\text{DST-I} &: \quad \mathbf{S}_{N-1}^I := \left( \frac{2}{N} \right)^{1/2} \left( \sin \frac{(j+1)(k+1)\pi}{N} \right)_{j,k=0}^{N-2} \in \mathbb{R}^{N-1,N-1}, \\
\text{DST-II} &: \quad \mathbf{S}_N^{II} := \left( \frac{2}{N} \right)^{1/2} \left( \varepsilon_{k+1}^N \sin \frac{(j+1)(2k+1)\pi}{2N} \right)_{j,k=0}^{N-1} \in \mathbb{R}^{N,N}, \\
\text{DST-III} &: \quad \mathbf{S}_N^{III} := (\mathbf{S}_N^{II})' \in \mathbb{R}^{N,N}, \\
\text{DST-IV} &: \quad \mathbf{S}_N^{IV} := \left( \frac{2}{N} \right)^{1/2} \left( \cos \frac{(2j+1)(2k+1)\pi}{4N} \right)_{j,k=0}^{N-1} \in \mathbb{R}^{N,N},
\end{aligned}$$

where  $\varepsilon_k^N := 1/\sqrt{2}$  ( $k = 0, N$ ) and  $\varepsilon_k^N := 1$  otherwise. We refer to the corresponding transforms as *trigonometric transforms*. It is well-known that the above matrices are orthogonal and that the vector multiplication with any of these matrices takes only  $O(N \log N)$  arithmetical operations. Fortunately, there exist implementations of algorithms for the vector multiplication with the above sine and cosine matrices, for example a C-implementation based on [3] and [19].

Moreover, we use the slightly modified DCT-I and DST-I matrices

$$\tilde{\mathbf{C}}_{N+1}^I := \left( (\varepsilon_k^N)^2 \cos \frac{jk\pi}{N} \right)_{j,k=0}^N, \quad \tilde{\mathbf{S}}_{N-1}^I := \left( \sin \frac{jk\pi}{N} \right)_{j,k=1}^{N-1}$$

and the slightly modified DCT-III and DST-III matrices

$$\tilde{\mathbf{C}}_N^{III} := \left( (\varepsilon_k^N)^2 \cos \frac{(2j+1)k\pi}{2N} \right)_{j,k=0}^{N-1}, \quad \tilde{\mathbf{S}}_N^{III} := \left( (\varepsilon_{k+1}^N)^2 \sin \frac{(2j+1)(k+1)\pi}{2N} \right)_{j,k=0}^{N-1}.$$

It holds that

$$\tilde{\mathbf{C}}_{N+1}^I \tilde{\mathbf{C}}_{N+1}^I = \frac{N}{2} \mathbf{I}_{N+1}. \quad (3.1)$$

**Theorem 3.1.** There exist the following relations between trigonometric transforms and Toeplitz matrices:

i) DCT-I and DST-I:

$$\begin{aligned}
\mathbf{R}'_N \mathbf{C}_{N+1}^I \mathbf{D} \mathbf{C}_{N+1}^I \mathbf{R}_N &= \frac{1}{2} \text{stoeplitz}(a_0, \dots, a_{N-2}) + \frac{1}{2} \text{shankel}(a_2, \dots, a_{N-2}, 0, 0), \\
\mathbf{S}_{N-1}^I \mathbf{R}'_N \mathbf{D} \mathbf{R}_N \mathbf{S}_{N-1}^I &= \frac{1}{2} \text{stoeplitz}(a_0, \dots, a_{N-2}) - \frac{1}{2} \text{shankel}(a_2, \dots, a_{N-2}, 0, 0), \\
\mathbf{R}'_N \mathbf{C}_{N+1}^I \tilde{\mathbf{D}} \mathbf{R}_N \mathbf{S}_{N-1}^I &= \frac{1}{2} \text{atoeplitz}(0, a_1, \dots, a_{N-2}) + \frac{1}{2} \text{ahankel}(a_2, \dots, a_{N-1}, 0), \\
\mathbf{S}_{N-1}^I \mathbf{R}'_N \tilde{\mathbf{D}} \mathbf{C}_{N+1}^I \mathbf{R}_N &= -\frac{1}{2} \text{atoeplitz}(0, a_1, \dots, a_{N-2}) + \frac{1}{2} \text{ahankel}(a_2, \dots, a_{N-1}, 0)
\end{aligned}$$

with

$$\begin{aligned}
\mathbf{D} &:= \text{diag}(d_0, \dots, d_N), \quad \tilde{\mathbf{D}} := \text{diag}(0, \tilde{d}_1, \dots, \tilde{d}_{N-1}, 0), \\
(d_0, \dots, d_N)' &:= \tilde{\mathbf{C}}_{N+1}^I (a_0, \dots, a_{N-2}, 0, 0)', \\
(\tilde{d}_1, \dots, \tilde{d}_{N-1})' &:= \tilde{\mathbf{S}}_{N-1}^I (a_1, \dots, a_{N-1})'.
\end{aligned}$$

ii) DCT-II and DST-II:

$$\begin{aligned} (\mathbf{C}_N^{II})' \mathbf{Z}'_{N,2} \mathbf{D} \mathbf{Z}_{N,2} \mathbf{C}_N^{II} &= \frac{1}{2} \text{stoeplitz}(a_0, \dots, a_{N-1}) + \frac{1}{2} \text{shankel}(a_1, \dots, a_{N-1}, 0), \\ (\mathbf{S}_N^{II})' \mathbf{Z}'_{N,1} \mathbf{D} \mathbf{Z}_{N,1} \mathbf{S}_N^{II} &= \frac{1}{2} \text{stoeplitz}(a_0, \dots, a_{N-1}) - \frac{1}{2} \text{shankel}(a_1, \dots, a_{N-1}, 0), \\ (\mathbf{C}_N^{II})' \mathbf{Z}'_{N,2} \tilde{\mathbf{D}} \mathbf{Z}_{N,1} \mathbf{S}_N^{II} &= \frac{1}{2} \text{atoeplitz}(0, a_1, \dots, a_{N-1}) + \frac{1}{2} \text{ahankel}(a_1, \dots, a_{N-1}, 0), \\ (\mathbf{S}_N^{II})' \mathbf{Z}'_{N,1} \tilde{\mathbf{D}} \mathbf{Z}_{N,2} \mathbf{C}_N^{II} &= -\frac{1}{2} \text{atoeplitz}(0, a_1, \dots, a_{N-1}) + \frac{1}{2} \text{ahankel}(a_1, \dots, a_{N-1}, 0) \end{aligned}$$

with

$$\begin{aligned} \mathbf{D} &:= \text{diag}(d_0, \dots, d_N), \quad \tilde{\mathbf{D}} := \text{diag}(0, \tilde{d}_1, \dots, \tilde{d}_{N-1}, 0), \\ (d_0, \dots, d_N)' &:= \tilde{\mathbf{C}}_{N+1}^I (a_0, \dots, a_{N-1}, 0)', \\ (\tilde{d}_1, \dots, \tilde{d}_{N-1})' &:= \tilde{\mathbf{S}}_{N-1}^I (a_1, \dots, a_{N-1})'. \end{aligned}$$

iii) DCT-IV and DST-IV:

$$\begin{aligned} \mathbf{C}_N^{IV} \mathbf{D} \mathbf{C}_N^{IV} &= \frac{1}{2} \text{stoeplitz}(a_0, \dots, a_{N-1}) + \frac{1}{2} \text{shankel}(a_1, \dots, a_{N-1}, 0), \\ \mathbf{S}_N^{IV} \mathbf{D} \mathbf{S}_N^{IV} &= \frac{1}{2} \text{stoeplitz}(a_0, \dots, a_{N-1}) - \frac{1}{2} \text{shankel}(a_1, \dots, a_{N-1}, 0), \\ \mathbf{C}_N^{IV} \tilde{\mathbf{D}} \mathbf{S}_N^{IV} &= -\frac{1}{2} \text{atoeplitz}(0, a_1, \dots, a_{N-1}) + \frac{1}{2} \text{shankel}(a_1, \dots, a_{N-1}, 0), \\ \mathbf{S}_N^{IV} \tilde{\mathbf{D}} \mathbf{C}_N^{IV} &= \frac{1}{2} \text{atoeplitz}(0, a_1, \dots, a_{N-1}) + \frac{1}{2} \text{shankel}(a_1, \dots, a_{N-1}, 0) \end{aligned}$$

with

$$\begin{aligned} \mathbf{D} &:= \text{diag}(d_0, \dots, d_{N-1}), \quad \tilde{\mathbf{D}} := \text{diag}(\tilde{d}_0, \dots, \tilde{d}_{N-1}) \\ (d_0, \dots, d_{N-1})' &:= \tilde{\mathbf{C}}_N^{III} (a_0, \dots, a_{N-1})', \\ (\tilde{d}_0, \dots, \tilde{d}_{N-1})' &:= \tilde{\mathbf{S}}_N^{III} (a_1, \dots, a_{N-1}, 0)'. \end{aligned}$$

Note that for fixed diagonal matrices  $\mathbf{D}, \tilde{\mathbf{D}}$ , the above decompositions into a Toeplitz and a Hankel matrix are not unique.

**Proof:** We restrict the proof to the DCT-II. To simplify the notation, we drop the index  $N$  and set  $\mathbf{C} := \mathbf{Z}_{N,2} \mathbf{C}_N^{II}$ ,  $\mathbf{S} := \mathbf{Z}_{N,1} \mathbf{S}_N^{II}$  and  $\mathbf{D} := \text{diag}(d_0, \dots, d_N)$ . Then by

$$\cos \alpha \cos \beta = \frac{1}{2} \cos(\alpha - \beta) + \frac{1}{2} \cos(\alpha + \beta),$$

the  $(u, v)$ -entry of the matrix  $\mathbf{C}' \mathbf{D} \mathbf{C}$  is

$$(\mathbf{C}' \mathbf{D} \mathbf{C})_{u,v} = \frac{1}{2} \frac{2}{N} \sum_{k=0}^{N-1} (\varepsilon_k^N)^2 d_k \cos \frac{(u-v)k\pi}{N} + \frac{1}{2} \frac{2}{N} \sum_{k=0}^{N-1} (\varepsilon_k^N)^2 d_k \cos \frac{(u+v+1)k\pi}{N},$$

or equivalently, since  $-(-1)^{u-v}d_N = -(-1)^{u+v+1}d_N$  for arbitrary  $d_N \in \mathbb{R}$ ,

$$(\mathbf{C}' \mathbf{D} \mathbf{C})_{u,v} = \frac{1}{2} \frac{2}{N} \sum_{k=0}^N (\varepsilon_k^N)^2 d_k \cos \frac{(u-v)k\pi}{N} + \frac{1}{2} \frac{2}{N} \sum_{k=0}^N (\varepsilon_k^N)^2 d_k \cos \frac{(u+v+1)k\pi}{N}.$$

Choosing  $d_N \in \mathbb{R}$  such that

$$\sum_{k=0}^N (\varepsilon_k^N)^2 d_k (-1)^k = 0,$$

we get by symmetry properties of cosine function that

$$\mathbf{C}' \mathbf{D} \mathbf{C} = \frac{1}{2} \text{stoeplitz}(a_0, \dots, a_{N-1}) + \frac{1}{2} \text{shankel}(a_1, \dots, a_{N-1}, 0),$$

where

$$(a_0, \dots, a_{N-1}, 0)' = \frac{2}{N} \tilde{\mathbf{C}}_{N+1}^I (d_0, \dots, d_N)',$$

i.e. by (3.1)

$$(d_0, \dots, d_N)' = \tilde{\mathbf{C}}_{N+1}^I (a_0, \dots, a_{N-1}, 0)'.$$

The other decomposition relations follow in a similar way by application of  $\sin \alpha \sin \beta = \frac{1}{2} \cos(\alpha - \beta) - \frac{1}{2} \cos(\alpha + \beta)$  and  $\sin \alpha \cos \beta = \frac{1}{2} \sin(\alpha - \beta) + \frac{1}{2} \sin(\alpha + \beta)$ . ■

Theorem 3.1 provides a new method for the fast multiplication of a real vector with a real nonsymmetric Toeplitz matrix that avoids the complex arithmetic which comes into the play if we exploit the usual FFT-based method for the fast vector – Toeplitz matrix multiplication.

**Corollary 3.2.** (Fast vector multiplication with nonsymmetric Toeplitz matrices)  
Let

$$\mathbf{T} = \mathbf{T}_N := (t_{j-k})_{j,k=0}^{N-1} = \text{toeplitz}((t_0, t_{-1}, \dots, t_{-(N-1)}), (t_0, t_1, \dots, t_{N-1}))$$

be given and let  $\mathbf{C} := \mathbf{Z}_{N,2} \mathbf{C}_N^{II}$ ,  $\mathbf{S} := \mathbf{Z}_{N,1} \mathbf{S}_N^{II}$ . Then

$$\mathbf{T} = \frac{1}{2} (\mathbf{T} + \mathbf{T}') + \frac{1}{2} (\mathbf{T} - \mathbf{T}') = \mathbf{C}' \mathbf{D} \mathbf{C} + \mathbf{S}' \mathbf{D} \mathbf{S} + \mathbf{C}' \tilde{\mathbf{D}} \mathbf{S} - \mathbf{S}' \tilde{\mathbf{D}} \mathbf{C},$$

where

$$\mathbf{D} := \text{diag}(d_0, \dots, d_N), \quad \tilde{\mathbf{D}} := \text{diag}(0, \tilde{d}_1, \dots, \tilde{d}_{N-1}, 0),$$

$$\begin{aligned} (d_0, \dots, d_N)' &:= \tilde{\mathbf{C}}_{N+1}^I (t_0, \frac{t_1 + t_{-1}}{2}, \dots, \frac{t_{N-1} + t_{-(N-1)}}{2}, 0)', \\ (\tilde{d}_1, \dots, \tilde{d}_{N-1})' &:= \tilde{\mathbf{S}}_{N-1}^I (\frac{t_{-1} - t_1}{2}, \dots, \frac{t_{-(N-1)} - t_{N-1}}{2})'. \end{aligned}$$

The vector multiplication with  $\mathbf{T}$  requires except of  $O(N)$  additions  
- one DCT-I and one DST-I to build  $\mathbf{D}$  and  $\tilde{\mathbf{D}}$  in a precomputation step,

- one DCT-II and one DST-II,
- four multiplications of vectors with diagonal matrices,
- one DCT-III and one DST-III of the vectors  $\mathbf{DC}\mathbf{x} + \tilde{\mathbf{DS}}\mathbf{x}$  and  $\mathbf{DS}\mathbf{x} - \tilde{\mathbf{DC}}\mathbf{x}$ , respectively, and takes therefore only  $O(N \log N)$  arithmetical operations.

Clearly, by Theorem 3.1, we can formulate similar algorithms for the fast multiplication of vectors with Toeplitz or Hankel matrices with respect to the other trigonometric transforms. Typewriting this paper, we got a ps-file of a paper of G. Heinig and K. Rost [14], which contains results in a similar direction as presented in this section.

## 4 Optimal trigonometric preconditioners

We are concerned with the solution of the system of linear equations

$$\mathbf{T}_N \mathbf{x}_N = \mathbf{b}_N$$

with a nonsingular nonsymmetric Toeplitz matrix  $\mathbf{T}_N \in \mathbb{R}^{N,N}$ . We intend to solve the normal equation

$$\mathbf{T}'_N \mathbf{T}_N \mathbf{x}_N = \mathbf{T}'_N \mathbf{b}_N \quad (4.1)$$

by the CG-method. In Section 7, we will see that with a good preconditioner at hand, this can be realized in a fast way. There are several requirements on a preconditioner  $\mathbf{M}_N$  of (4.1) resulting from the construction and the convergence behaviour of the CG-method as well as from the fact that the vector multiplication with  $\mathbf{T}_N$  requires only  $O(N \log N)$  arithmetical operations. Therefore, we are looking for a preconditioner with the following properties:

(P1)  $\mathbf{M}_N$  is symmetric and positive definite such that the bilinear form

$$(\mathbf{x}_N, \mathbf{y}_N)_{\mathbf{M}_N} := \mathbf{x}'_N \mathbf{M}_N \mathbf{y}_N$$

arising in the left preconditioned CG-method is symmetric and positive definite, too.

(P2) The spectrum of  $\mathbf{M}_N^{-1} \mathbf{T}'_N \mathbf{T}_N$  is clustered at 1.

(P3) The vector multiplication with  $\mathbf{M}_N$  can be computed with  $O(N \log N)$  arithmetical operations.

(P4) The construction of  $\mathbf{M}_N$  takes only  $O(N \log N)$  arithmetical operations.

Having property (P3) in mind, a straightforward idea consists in choosing  $\mathbf{M}_N$  from an algebra

$$\mathcal{A}_{\mathbf{O}_N} := \{\mathbf{O}'_N (\text{diag } \mathbf{d}) \mathbf{O}_N : \mathbf{d} \in \mathbb{R}^N\} \quad (4.2)$$

of matrices which are diagonalizable by some orthogonal matrix  $\mathbf{O}_N$ , where  $\mathbf{O}_N$  has the additional property that its vector multiplication requires only  $O(N \log N)$  arithmetical operations. As orthogonal matrices, we will use the trigonometric matrices of the previous section which are closely related to the Fourier matrix  $\mathbf{F}_N$ , but have the advantage of purely real entries. Moreover, if we choose  $\mathbf{M}_N \in \mathcal{A}_{\mathbf{O}_N}$  as so-called optimal preconditioner of  $\mathbf{T}'_N \mathbf{T}_N$ , then we will see that under certain assumptions on  $\mathbf{T}_N$ , the properties (P1), (P2) and (P4) are also fulfilled.

For  $\mathbf{A}_N \in \mathbb{R}^{N,N}$ , the matrix  $\mathbf{M}_N(\mathbf{A}_N)$  is called an *optimal preconditioner* of  $\mathbf{A}_N$  in  $\mathcal{A}_{O_N}$  [5] if

$$\|\mathbf{M}_N(\mathbf{A}_N) - \mathbf{A}_N\|_F = \min\{\|\mathbf{B}_N - \mathbf{A}_N\|_F : \mathbf{B}_N \in \mathcal{A}_{O_N}\}. \quad (4.3)$$

If  $\mathbf{O}_N$  is one of the orthogonal matrices which correspond to the DST-I, DST-II, DCT-II, DST-IV or DCT-IV, respectively, then  $\mathbf{M}_N$  is said to be an *optimal trigonometric preconditioner* of  $\mathbf{A}_N$ . The choice of the Frobenius norm in definition (4.3) results from the fact that the Frobenius norm is induced by an inner product of  $\mathbb{R}^{N,N}$

$$\langle \mathbf{A}_N, \mathbf{B}_N \rangle := \text{tr}(\mathbf{A}'_N \mathbf{B}_N) = \sum_{j,k=0}^{N-1} a_{j,k} b_{j,k}. \quad (4.4)$$

In particular, it holds that

$$\|\mathbf{O}_N \mathbf{A}_N \mathbf{O}'_N\|_F^2 = \text{tr}(\mathbf{O}_N \mathbf{A}'_N \mathbf{O}'_N \mathbf{O}_N \mathbf{A}_N \mathbf{O}'_N) = \text{tr}(\mathbf{A}'_N \mathbf{A}_N) = \|\mathbf{A}_N\|_F^2. \quad (4.5)$$

The following lemma describes the optimal preconditioner of  $\mathbf{T}'_N \mathbf{T}_N$  in two different ways.

**Lemma 4.1.** Let  $\mathbf{A}_N \in \mathbb{R}^{N,N}$  and let  $\mathcal{A}_{O_N}$  be defined by (4.2) with respect to some orthogonal matrix  $\mathbf{O}_N$ . Then the optimal preconditioner of  $\mathbf{A}_N$  is given by

$$\mathbf{M}_N(\mathbf{A}_N) = \mathbf{O}'_N \delta(\mathbf{O}_N \mathbf{A}_N \mathbf{O}'_N) \mathbf{O}_N. \quad (4.6)$$

If  $\{\mathbf{B}_0, \dots, \mathbf{B}_{N-1}\}$  denotes a basis of  $\mathcal{A}_{O_N}$ , then an alternative description of  $\mathbf{M}_N(\mathbf{A}_N)$  reads as

$$\mathbf{M}_N(\mathbf{A}_N) := \sum_{k=0}^{N-1} \alpha_k \mathbf{B}_k, \quad (4.7)$$

where the coefficient vector  $\boldsymbol{\alpha} := (\alpha_0, \dots, \alpha_{N-1})'$  is determined by

$$\mathbf{G}\boldsymbol{\alpha} = \boldsymbol{\beta}, \quad \mathbf{G} := (\langle \mathbf{B}_k, \mathbf{B}_j \rangle)_{j,k=0}^{N-1}, \quad \boldsymbol{\beta} := (\langle \mathbf{A}_N, \mathbf{B}_j \rangle)_{j=0}^{N-1}.$$

**Proof:** 1. By (4.5), it follows for  $\mathbf{M}_N := \mathbf{O}'_N (\text{diag } \mathbf{d}) \mathbf{O}_N \in \mathcal{A}_{O_N}$  that

$$\|\mathbf{M}_N - \mathbf{A}_N\|_F = \|\text{diag } \mathbf{d} - \mathbf{O}_N \mathbf{A}_N \mathbf{O}'_N\|_F$$

which implies (4.6) by the definition of the optimal preconditioner.

2. The computation of the optimal preconditioner of  $\mathbf{A}_N$  in  $\mathcal{A}_{O_N}$  is equivalent with the computation of the element of best approximation of  $\mathbf{A}_N$  in the linear subspace  $\mathcal{A}_{O_N}$  of the Hilbert space  $\mathbb{R}^{N,N}$  equipped with the inner product (4.4). This can be done by the Galerkin approach (4.7). ■

Now it is easy to verify that an optimal preconditioner of  $\mathbf{T}'_N \mathbf{T}_N$  satisfies property (P1).

**Corollary 4.2.** Let  $\mathbf{A}_N \in R^{N,N}$  be a symmetric positive definite matrix and let  $\mathcal{A}_{O_N}$  be defined by (4.2) with respect to some orthogonal matrix  $\mathbf{O}_N$ . Then the optimal preconditioner  $\mathbf{M}_N = \mathbf{M}_N(\mathbf{A}_N)$  is also symmetric and positive definite.

**Proof:** The symmetry of  $\mathbf{M}_N$  follows by definition of  $\mathcal{A}_{O_N}$ .

By (4.6), the eigenvalues of  $\mathbf{M}_N$  are given by the diagonal entries  $(\mathbf{O}_N \mathbf{A}_N \mathbf{O}'_N)_{k,k}$  ( $k = 0, \dots, N-1$ ) of  $\mathbf{O}_N \mathbf{A}_N \mathbf{O}'_N$ . Since  $\mathbf{A}_N$  is positive definite, it holds that

$$\begin{aligned} 0 &< \min \left\{ \frac{\mathbf{x}'_N \mathbf{A}_N \mathbf{x}_N}{\mathbf{x}'_N \mathbf{x}_N} : \mathbf{x}_N \neq \mathbf{0}_N \right\} = \min \left\{ \frac{\mathbf{y}'_N \mathbf{O}_N \mathbf{A}_N \mathbf{O}'_N \mathbf{y}_N}{\mathbf{y}'_N \mathbf{y}_N} : \mathbf{y}_N \neq \mathbf{0}_N \right\} \\ &\leq \mathbf{e}'_k (\mathbf{O}_N \mathbf{A}_N \mathbf{O}'_N) \mathbf{e}_k = (\mathbf{O}_N \mathbf{A}_N \mathbf{O}'_N)_{k,k} \quad (k = 0, \dots, N-1), \end{aligned}$$

and we are done. ■

In connection with (4.6) and Corollary 4.2 see also [18].

## 5 Clusters of eigenvalues

Let  $C_{2\pi}$  denote the Banach space of  $2\pi$ -periodic complex-valued functions equipped with the usual norm  $\|\cdot\|_\infty$ . In this section, we are only interested in functions  $f = f_R + i f_I \in C_{2\pi}$  with real Fourier coefficients

$$t_k := \int_{-\pi}^{\pi} f(x) e^{-ikx} dx \quad (k \in \mathbb{Z}),$$

where we suppose that the real and the imaginary part of  $f$

$$\begin{aligned} f_R &= t_0 + \sum_{k=1}^{\infty} (t_k + t_{-k}) \cos kx, \\ f_I &= \sum_{k=1}^{\infty} (t_k - t_{-k}) \sin kx \end{aligned}$$

do not vanish, respectively. Moreover, we assume that  $\{t_k\}_{k \in \mathbb{Z}} \in l_1$ , such that  $f$  belongs to the Wiener class. Consider the  $N$ -th Toeplitz matrix corresponding to the generating function  $f$

$$\mathbf{T}_N := \text{toeplitz}((t_0, t_{-1}, \dots, t_{-(N-1)}), (t_0, t_1, \dots, t_{N-1})).$$

It is well-known that the singular values of  $T_N$  are distributed as  $|f|$  [20]. Note that the above result was extended to functions  $f \in L^2_{2\pi} \supset C_{2\pi}$  in [22].

The following definition is due to E.E. Tyrtyshnikov [22]. Let  $\{\sigma_k^N\}_{k=1}^N$  be a sequence of real numbers and let  $\gamma_N(\varepsilon)$  denote the number of those among  $\sigma_k^N$  ( $k = 1, \dots, N$ ) which are outside the  $\varepsilon$ -ball centered at  $p$ . If  $\gamma_N(\varepsilon) < K(\varepsilon)$ , where  $K(\varepsilon)$  is independent of  $N$ , then  $p$  is called a *proper cluster*. In this sense, we say that the values  $\sigma_k^N$  are clustered at  $p$ .

In the following, we restrict our attention to preconditioners of  $\mathcal{A}_{C_N^{II}}$ . By Theorem 3.1, the approach for the preconditioners with respect to the DST-I, the DST-II, the

DST-IV and the DCT-IV follows the same lines. Let  $\mathbf{C} = \mathbf{C}_N := \mathbf{Z}_{N,2}\mathbf{C}_N^{II}$  and  $\mathbf{S} = \mathbf{S}_N := \mathbf{Z}_{N,1}\mathbf{S}_N^{II}$ . Regarding Theorem 3.1 ii), we associate with the sequence  $\{\mathbf{T}_N\}_{N \in \mathbb{N}}$  of Toeplitz matrices a sequence  $\{\mathbf{H}_N\}_{N \in \mathbb{N}}$  of Hankel matrices

$$\mathbf{H}_N := \text{hankel}((t_{-1}, \dots, t_{-(N-1)}, 0), (t_1, \dots, t_{N-1}, 0)).$$

Let  $\mathbf{M}_N$  denote the optimal preconditioner of  $\mathbf{T}'_N \mathbf{T}_N$  with respect to the DCT-II, i.e.

$$\mathbf{M}_N := (\mathbf{C}_N^{II})' \delta(\mathbf{C}_N^{II} \mathbf{T}'_N \mathbf{T}_N (\mathbf{C}_N^{II})') \mathbf{C}_N^{II}. \quad (5.1)$$

In this section, we prove that under certain assumptions on  $\mathbf{T}_N$ , the eigenvalues of the preconditioned matrix  $\mathbf{M}_N^{-1} \mathbf{T}'_N \mathbf{T}_N$  are clustered at 1. We follow the lines of R.H. Chan. First, we show that for all  $\varepsilon > 0$  and  $N$  sufficiently large, the matrix  $\mathbf{T}'_N \mathbf{T}_N - \mathbf{M}_N$  splits into a matrix of low rank independent of  $N$  and a matrix with spectral norm smaller than  $\varepsilon$ . Then we apply Cauchy's interlace theorem to verify that the eigenvalues of

$$\mathbf{M}_N^{-1}(\mathbf{T}'_N \mathbf{T}_N - \mathbf{M}_N) = \mathbf{M}_N^{-1} \mathbf{T}'_N \mathbf{T}_N - \mathbf{I}_N$$

are clustered at 0. Again, we drop the index  $N$ , if the dimension of the matrices follows from the context.

In preparation of Theorem 5.3, we provide the following two lemmata.

**Lemma 5.1.** Let  $a = \{a_k\}_{k=0}^\infty \in l_1$  and  $b = \{b_k\}_{k=0}^\infty \in l_1$ . Then, for all  $\varepsilon > 0$ , there exists  $m = m(\varepsilon)$  such that for all  $N \geq 2m$  the Hankel matrix

$$\mathbf{H} := \text{hankel}((a_0, a_1, \dots, a_{N-1} + b_{N-1}), (b_0, b_1, \dots, a_{N-1} + b_{N-1}))$$

splits as  $\mathbf{H} = \mathbf{V}_H + \mathbf{W}_H$ , where

$$\mathbf{W}_H := \text{hankel}((a_0, \dots, a_{m-1}, \mathbf{0}_{N-m}), (b_0, \dots, b_{m-1}, \mathbf{0}_{N-m}))$$

is a matrix of rank  $\leq 2m$  and where  $\mathbf{V}_H := \mathbf{H} - \mathbf{W}_H$  satisfies  $\|\mathbf{V}_H\|_2 < \varepsilon$ .

**Proof:** Since  $a, b \in l_1$ , there exists for all  $\varepsilon > 0$  an integer  $m = m(\varepsilon)$  such that

$$\sum_{k=m}^{\infty} |a_k| < \varepsilon/2, \quad \sum_{k=m}^{\infty} |b_k| < \varepsilon/2.$$

Now the assertion follows from  $\|\mathbf{V}_H\|_2 \leq \|\mathbf{V}_H\|_1 < \varepsilon/2 + \varepsilon/2$ . ■

**Lemma 5.2.** Let  $t = \{t_k\}_{k \in \mathbb{Z}} \in l_1$  with  $\|t\|_1 = \tau$  be given. Further let  $\mathbf{T} = \mathbf{T}_N$  and  $\mathbf{H} = \mathbf{H}_N$  be the corresponding  $N$ -th Toeplitz matrix and  $N$ -th Hankel matrix, respectively. Then, for all  $\varepsilon > 0$ , there exists  $m = m(\varepsilon)$  such that for all  $N \geq 4m$

$$\mathbf{H}\mathbf{T} + \mathbf{T}'\mathbf{H} + \mathbf{H}^2 = \mathbf{V} + \mathbf{W}, \quad (5.2)$$

where  $\|\mathbf{V}\|_2 < \varepsilon$  and

$$\mathbf{W} := (w_{j,k})_{j,k=0}^{N-1} \text{ with } w_{j,k} = 0 \text{ for } 2m \leq j+k \leq 2N-2-m. \quad (5.3)$$

**Proof:** By construction, it holds that  $\|\mathbf{T}\|_2 \leq \tau$ ,  $\|\mathbf{H}\|_2 \leq \tau$ . Since  $t \in l_1$ , there exists  $m = m(\varepsilon)$  such that

$$\sum_{|k|=m+1}^{\infty} |t_k| < \frac{\varepsilon}{6\tau}. \quad (5.4)$$

Then

$$\mathbf{T} = \mathbf{T}_\varepsilon + \mathbf{T}_B, \quad \mathbf{H} = \mathbf{H}_\varepsilon + \mathbf{H}_B \quad (5.5)$$

with

$$\begin{aligned} \mathbf{T}_\varepsilon &:= \text{toeplitz}((\mathbf{0}_{m+1}, t_{-(m+1)}, \dots, t_{-(N-1)}), (\mathbf{0}_{m+1}, t_{m+1}, \dots, t_{N-1})), \\ \mathbf{T}_B &:= \text{toeplitz}((t_0, \dots, t_{-m}, \mathbf{0}_{N-m-1}), (t_0, \dots, t_m, \mathbf{0}_{N-m-1})), \\ \mathbf{H}_\varepsilon &:= \text{hankel}((\mathbf{0}_m, t_{-(m+1)}, \dots, t_{-(N-2)}, 0), (\mathbf{0}_m, t_{m+1}, \dots, t_{N-2}, 0)), \\ \mathbf{H}_B &:= \text{hankel}((t_{-1}, \dots, t_{-m}, \mathbf{0}_{N-m}), (t_1, \dots, t_m, \mathbf{0}_{N-m})), \end{aligned}$$

where we obtain by (5.4) that  $\|\mathbf{T}_\varepsilon\|_2 \leq \frac{\varepsilon}{6\tau}$ ,  $\|\mathbf{H}_\varepsilon\|_2 \leq \frac{\varepsilon}{6\tau}$ . Substituting (5.5) in (5.2), we obtain the desired decomposition

$$\begin{aligned} \mathbf{HT} + \mathbf{T}'\mathbf{H} + \mathbf{H}^2 &= (\mathbf{H}_\varepsilon + \mathbf{H}_B)(\mathbf{T}_\varepsilon + \mathbf{T}_B) + (\mathbf{T}'_\varepsilon + \mathbf{T}'_B)(\mathbf{H}_\varepsilon + \mathbf{H}_B) + (\mathbf{H}_\varepsilon + \mathbf{H}_B)^2 \\ &= (\mathbf{HT}_\varepsilon + \mathbf{T}'_\varepsilon\mathbf{H} + \mathbf{H}_\varepsilon\mathbf{T}_B + \mathbf{T}'_B\mathbf{H}_\varepsilon + \mathbf{H}_\varepsilon\mathbf{H} + \mathbf{H}_B\mathbf{H}_\varepsilon) \\ &\quad + (\mathbf{T}'_B\mathbf{H}_B + \mathbf{H}_B\mathbf{T}_B + \mathbf{H}_B^2). \end{aligned}$$

■

**Theorem 5.3.** Let  $t = \{t_k\}_{k \in \mathbb{Z}} \in l_1$  with  $\|t\|_1 = \tau$  be given and let  $\mathbf{T} = \mathbf{T}_N$  be the corresponding  $N$ -th Toeplitz matrix. Then, for all  $\varepsilon > 0$ , there exists  $m = m(\varepsilon)$  such that for all  $N \geq 2m$

$$\mathbf{T}'\mathbf{T} = \mathbf{C}'\mathbf{D}\mathbf{C} + \mathbf{V} + \mathbf{W},$$

where  $\mathbf{D}$  denotes some diagonal matrix,

$$\mathbf{W} := (w_{j,k})_{j,k=0}^{N-1} \text{ with } w_{j,k} = 0 \text{ for } m \leq j+k \leq 2N-2-m,$$

and where  $\|\mathbf{V}\|_2 < \varepsilon$ .

**Proof:** Let  $\mathbf{H} = \mathbf{H}_N$  denote the  $N$ -Hankel matrix associated with  $t$ . Then it follows by Theorem 3.1 ii) that

$$\begin{aligned} \mathbf{T}'\mathbf{T} &= (\mathbf{T}' + \mathbf{H} - \mathbf{H})(\mathbf{T} + \mathbf{H} - \mathbf{H}) \\ &= (\mathbf{C}'\mathbf{D}_a\mathbf{C} + \mathbf{S}'\tilde{\mathbf{D}}_b\mathbf{C} - \mathbf{H})(\mathbf{C}'\mathbf{D}_a\mathbf{C} + \mathbf{C}'\tilde{\mathbf{D}}_b\mathbf{S} - \mathbf{H}), \end{aligned}$$

where

$$\begin{aligned}\mathbf{a} = \mathbf{a}_{N+1} &:= (2t_0, t_{-1} + t_1, \dots, t_{-(N-1)} + t_{N-1}, 0)' \in \mathbb{R}^{N+1}, \\ \mathbf{b} = \mathbf{b}_{N-1} &:= (t_{-1} - t_1, \dots, t_{-(N-1)} - t_{N-1})' \in \mathbb{R}^{N-1}\end{aligned}$$

and

$$\begin{aligned}\mathbf{D}_a &= \text{diag}(d_0, \dots, d_{N-1}, d_N) \quad , \quad (d_0, \dots, d_N)' := \tilde{\mathbf{C}}_{N+1}^I \mathbf{a}, \\ \tilde{\mathbf{D}}_b &= \text{diag}(0, \tilde{d}_1, \dots, \tilde{d}_{N-1}, 0) \quad , \quad (\tilde{d}_0, \dots, \tilde{d}_{N-1})' := \tilde{\mathbf{S}}_{N-1}^I \mathbf{b}.\end{aligned}$$

By  $\mathbf{C}_N^{II}(\mathbf{C}_N^{II})' = \mathbf{I}_N$ , we further obtain that

$$\mathbf{T}'\mathbf{T} = \mathbf{C}'\mathbf{D}_a^2\mathbf{C} + \mathbf{S}'\tilde{\mathbf{D}}_b^2\mathbf{S} + \mathbf{C}'\mathbf{D}_a\tilde{\mathbf{D}}_b\mathbf{S} + \mathbf{S}'\tilde{\mathbf{D}}_b\mathbf{D}_a\mathbf{C} - \mathbf{H}(\mathbf{T} + \mathbf{H}) - (\mathbf{T}' + \mathbf{H})\mathbf{H} + \mathbf{H}^2$$

and by Theorem 3.1 ii) that

$$\mathbf{T}'\mathbf{T} = \mathbf{C}'\mathbf{D}_a^2\mathbf{C} + \mathbf{C}'\tilde{\mathbf{D}}_b^2\mathbf{C} - \mathbf{H}_{\tilde{\mathbf{D}}_b^2} + \tilde{\mathbf{H}}_{\tilde{\mathbf{D}}_b\mathbf{D}_a} - (\mathbf{H}\mathbf{T} + \mathbf{T}'\mathbf{H} + \mathbf{H}^2) \quad (5.6)$$

with

$$\begin{aligned}\mathbf{H}_{\tilde{\mathbf{D}}_b^2} &:= \text{shankel}(u_1, \dots, u_N), \\ (u_0, \dots, u_N)' &:= \frac{2}{N} \tilde{\mathbf{C}}_{N+1}^I(0, \tilde{d}_1^2, \dots, \tilde{d}_{N-1}^2, 0)', \\ \tilde{\mathbf{H}}_{\tilde{\mathbf{D}}_b\mathbf{D}_a} &:= \text{ahankel}(v_1, \dots, v_{N-1}, 0), \\ (v_1, \dots, v_{N-1})' &:= \frac{2}{N} \tilde{\mathbf{S}}_{N-1}^I(d_1\tilde{d}_1, \dots, d_{N-1}\tilde{d}_{N-1})'.\end{aligned}$$

Set

$$\hat{\mathbf{H}} := \mathbf{H}_{\tilde{\mathbf{D}}_b^2} - \tilde{\mathbf{H}}_{\tilde{\mathbf{D}}_b\mathbf{D}_a} = \text{hankel}((w_1, \dots, w_N), (w_{-1}, \dots, w_{-N})) \quad (5.7)$$

with  $w_N = w_{-N} = u_N$ ,  $w_k = u_k - v_k$ ,  $w_{-k} = u_k + v_k$  ( $k = 1, \dots, N-1$ ). Then we have for all  $N \in \mathbb{N}$  that

$$\begin{aligned}w_N &= \frac{2}{N} \sum_{k=1}^{N-1} (-1)^k \tilde{d}_k^2, \\ |w_N| &\leq \frac{2}{N} \|\tilde{\mathbf{S}}_{N-1}^I \mathbf{b}\|_2^2 \leq \|\mathbf{b}\|_2^2 \leq \|\mathbf{b}\|_1^2 \leq \tau^2.\end{aligned} \quad (5.8)$$

Moreover, we get by Theorem 3.1 i) for the first row  $\mathbf{w} := (w_1, \dots, w_{N-1})'$  of  $\hat{\mathbf{H}}$  that

$$\begin{aligned}\mathbf{w} &= \frac{2}{N} \mathbf{R}'_N \tilde{\mathbf{C}}_{N+1}^I \text{diag}(0, \tilde{d}_1, \dots, \tilde{d}_{N-1}, 0) \mathbf{R}_N \tilde{\mathbf{S}}_{N-1}^I \mathbf{b} \\ &\quad - \frac{2}{N} \tilde{\mathbf{S}}_{N-1}^I \mathbf{R}'_N \text{diag}(d_0, \dots, d_N) \mathbf{R}_N \tilde{\mathbf{S}}_{N-1}^I \mathbf{b} \\ &= \text{toeplitz}((-t_0, -t_1, \dots, -t_{N-2}), (-t_0, -t_{-1}, \dots, -t_{-(N-2)})) \mathbf{b} \\ &\quad + \text{hankel}((t_{-2}, \dots, t_{-(N-2)}, b_{N-1}, 0), (t_2, \dots, t_{N-2}, -b_{N-1}, 0)) \mathbf{b}.\end{aligned}$$

Thus, it holds for all  $N \in \mathbb{N}$  that

$$\|\mathbf{w}\|_1 \leq 2\tau^2.$$

Similarly, we conclude that

$$\|(w_{-1}, \dots, w_{-(N-1)})'\|_1 \leq \tau^2.$$

Together with (5.8), we see that  $\hat{\mathbf{H}}$  satisfies the assertion of Lemma 5.1. Hence, for fixed  $\varepsilon > 0$ , there exists  $\hat{m} = \hat{m}(\varepsilon) > 0$  such that for all  $N \geq 2\hat{m}$

$$\hat{\mathbf{H}} = \hat{\mathbf{V}} + \hat{\mathbf{W}} \quad (5.9)$$

with  $\|\hat{\mathbf{V}}\|_2 \leq \varepsilon/2$ ,  $\hat{\mathbf{W}} = \text{hankel}((w_1, \dots, w_m, \mathbf{0}_{N-m}), (w_{-1}, \dots, w_{-m}, \mathbf{0}_{N-m}))$ . Furthermore, by Lemma 5.2, there exists  $\tilde{m} = \tilde{m}(\varepsilon) > 0$  such that for  $N$  sufficiently large

$$(\mathbf{H}\mathbf{T} + \mathbf{T}'\mathbf{H} + \mathbf{H}^2) = \tilde{\mathbf{V}} + \tilde{\mathbf{W}} \quad (5.10)$$

with  $\|\tilde{\mathbf{V}}\|_2 \leq \varepsilon/2$  and with a low rank matrix  $\tilde{\mathbf{W}}$  of the form (5.3). Applying (5.9) and (5.10) in (5.6), we obtain the assertion

$$\mathbf{T}'\mathbf{T} = \mathbf{C}'(\mathbf{D}_a^2 + \tilde{\mathbf{D}}_b^2)\mathbf{C} + (\tilde{\mathbf{V}} - \hat{\mathbf{V}}) + (\tilde{\mathbf{W}}) - \hat{\mathbf{W}}$$

with  $m := \max\{\hat{m}, 2\tilde{m}\}$ . ■

**Lemma 5.4.** For  $m > 0$  and  $N > 2m$ , let  $\mathbf{V} \in \mathbb{R}^{N,N}$  with  $\|\mathbf{V}\|_2 < \varepsilon/2$  and

$$\mathbf{W} := (w_{j,k})_{j,k=0}^{N-1} \text{ with } w_{j,k} = 0 \text{ for } m \leq j+k \leq 2N-2-m,$$

be given. Set  $\omega := \sum_{j,k=0}^{N-1} |w_{j,k}|$ . Then it holds for  $N > 4\omega/\varepsilon$  that

$$\|\delta(\mathbf{C}_N^{II}(\mathbf{V} + \mathbf{W})(\mathbf{C}_N^{II})')\|_2 < \varepsilon.$$

Note that for fixed  $m$ , the value  $\omega$  does not depend on  $N$ .

**Proof:** On the one hand, we obtain that

$$\|\delta(\mathbf{C}_N^{II}\mathbf{V}(\mathbf{C}_N^{II})')\|_2 \leq \|\mathbf{C}_N^{II}\mathbf{V}(\mathbf{C}_N^{II})'\|_2 = \|\mathbf{V}\|_2 < \varepsilon/2,$$

and on the other hand that

$$\begin{aligned} \|\delta(\mathbf{C}_N^{II}\mathbf{W}(\mathbf{C}_N^{II})')\|_2 &\leq \max_{n=0,\dots,N-1} \left| \frac{2}{N} (\varepsilon_n^N)^2 \sum_{j,k=0}^{N-1} w_{j,k} \cos \frac{n(2j+1)\pi}{2N} \cos \frac{n(2k+1)\pi}{2N} \right| \\ &\leq 2\omega/N < \varepsilon/2 \quad (N > 4\omega/\varepsilon). \end{aligned}$$

Now summation implies the assertion. ■

**Theorem 5.5.** Let  $t = \{t_k\}_{k \in \mathbb{Z}} \in l_1$  with  $\|t\|_1 = \tau$  be given and let  $\mathbf{T}_N$  be the corresponding  $N$ -th Toeplitz matrix. Moreover, assume that the singular values of  $\mathbf{T}_N$  are larger than  $\gamma > 0$  for all  $N \in \mathbb{N}$ . Let  $\mathbf{M}_N = \mathbf{M}_N(\mathbf{T}'_N \mathbf{T}_N)$  denote the optimal preconditioner of  $\mathbf{T}'_N \mathbf{T}_N$  in  $\mathcal{A}_{C_N^{II}}$ . Then the eigenvalues of  $\mathbf{M}_N^{-1} \mathbf{T}'_N \mathbf{T}_N$  are clustered at 1.

**Proof:** By Corollary 4.2, it holds that  $\|\mathbf{M}_N^{-1}\|_2 < 1/\gamma$  for all  $N \in \mathbb{N}$ . Let  $\varepsilon > 0$  be fixed. Then we obtain by (5.1), Theorem 5.3 and Lemma 5.4 that for  $N$  sufficiently large, there exists  $M \in \mathbb{N}$  independent of  $N$  such that

$$\mathbf{T}'_N \mathbf{T}_N - \mathbf{M}_N = \mathbf{V}_N + \mathbf{W}_N - (\mathbf{C}_N^{II})' \delta(\mathbf{C}_N^{II} (\mathbf{V}_N + \mathbf{W}_N) (\mathbf{C}_N^{II})') \mathbf{C}_N^{II} = \mathbf{U}_N + \mathbf{W}_N,$$

where  $\mathbf{W}_N$  is of low rank  $M$  and where  $\|\mathbf{U}_N\|_2 < \varepsilon\gamma$ . Now

$$\begin{aligned} \mathbf{M}_N^{-1/2} \mathbf{T}'_N \mathbf{T}_N \mathbf{M}_N^{-1/2} - \mathbf{I}_N &= \mathbf{M}_N^{-1/2} \mathbf{U}_N \mathbf{M}_N^{-1/2} + \mathbf{M}_N^{-1/2} \tilde{\mathbf{W}}_N \mathbf{M}_N^{-1/2} \\ &= \tilde{\mathbf{V}}_N + \tilde{\mathbf{W}}_N, \end{aligned}$$

where  $\|\tilde{\mathbf{V}}_N\|_2 \leq \varepsilon$  and where the rank of  $\tilde{\mathbf{W}}_N$  is at most  $M$ . Since  $\tilde{\mathbf{V}}_N$  and  $\tilde{\mathbf{W}}_N$  are symmetric matrices, we can apply Cauchy's interlace theorem [25], which implies that for  $N$  sufficiently large, at most  $M$  eigenvalues of  $\tilde{\mathbf{V}}_N + \tilde{\mathbf{W}}_N$  have absolute value greater than  $\varepsilon$ . Now the assertion follows since  $\mathbf{M}_N^{-1/2} \mathbf{T}'_N \mathbf{T}_N \mathbf{M}_N^{-1/2}$  and  $\mathbf{M}_N^{-1} \mathbf{T}'_N \mathbf{T}_N$  possess the same eigenvalues. ■

**Remark.** Under the above assumptions on  $\mathbf{T}_N$  it was proved that the eigenvalues of  $(\mathbf{M}_N^{-1} \mathbf{T}_N)^*(\mathbf{M}_N^{-1} \mathbf{T}_N)$ , where  $\mathbf{M}_N = \mathbf{C}_N(\mathbf{T}_N)$  denotes the optimal circulant preconditioner of  $\mathbf{T}_N$ , are clustered at 1 [8, 22]. In general, the eigenvalues of  $(\mathbf{M}_N^{-1} \mathbf{T}_N)'(\mathbf{M}_N^{-1} \mathbf{T}_N)$  are not clustered at 1, if  $\mathbf{M}_N$  is the optimal trigonometric preconditioner of  $\mathbf{T}_N$ . If  $\mathbf{T}_N = (t_{j-k})_{j,k=0}^{N-1}$  with  $t_0 = 1$  and  $t_k = -t_{-k}$  ( $k = 1, \dots, N-1$ ), then the optimal trigonometric preconditioners of  $\mathbf{T}_N$  are  $\mathbf{M}_N = \mathbf{I}_N$ , i.e. we have no preconditioning. If, for example,  $t_{-1} = -t_1 = 2$  and  $t_k = 0$  ( $|k| > 1$ ), then the matrices  $\mathbf{T}_{2N+1}$  ( $N \in \mathbb{N}$ ) satisfy the assumption of Theorem 5.5. However, the eigenvalues of  $\mathbf{T}'_{2N+1} \mathbf{T}_{2N+1}$  are given by  $9 - 8 \cos(j\pi/(N+1))$  ( $j = 0, \dots, N$ ).

## 6 Construction of optimal preconditioners of $\mathbf{T}' \mathbf{T}$

In this section, we explain how optimal trigonometric preconditioners of  $\mathbf{T}'_N \mathbf{T}_N$  can be constructed with  $O(N \log N)$  arithmetical operations. In contrast to the construction of optimal trigonometric preconditioners of  $\mathbf{T}_N$ , we are confronted with the fact that  $\mathbf{T}'_N \mathbf{T}_N$  is not a Toeplitz matrix. Again, we consider only  $\mathbf{C}_N^{II}$ -preconditioners. The approach for the DST-I, DST-II, DST-IV and DCT-IV follows the same lines. For the construction of the optimal preconditioner  $\mathbf{M}_N = \mathbf{M}_N(\mathbf{T}'_N \mathbf{T}_N)$  we use the representation (4.7) of  $\mathbf{M}_N$  with the basis  $\{\mathbf{B}_k^{II} : k = 0, \dots, N-1\}$  of  $\mathcal{A}_{C_N^{II}}$  [2], [13]:

$$\mathbf{B}_k^{II} := (\mathbf{C}_N^{II})' \operatorname{diag}(U_k(c_l))_{l=0}^{N-1} \mathbf{C}_N^{II},$$

where  $c_l := \cos \frac{l\pi}{N}$  and where  $U_k$  denotes the  $k$ -th *Chebyshev polynomial of second kind*

$$U_k(x) := \sin((k+1)\arccos x)/\sin(\arccos x) \quad (x \in (-1, 1)).$$

Moreover, we apply that  $\{\mathbf{B}_k^I : k = 0, \dots, N-1\}$  with

$$\begin{aligned} \mathbf{B}_k^I &:= (\mathbf{C}_N^{II})' \operatorname{diag}(T_k(c_l))_{l=0}^{N-1} \mathbf{C}_N^{II} \\ &= \text{stoeplitz } \mathbf{e}'_k + \text{shankel } \mathbf{e}'_{k-1}, \quad (k = 1, \dots, N-1) \end{aligned} \quad (6.1)$$

$\mathbf{e}_{-1} := \mathbf{0}_N$ , and with the  $k$ -th *Chebyshev polynomial of first kind*

$$T_k(x) := \cos(k \arccos x) \quad (x \in [-1, 1]),$$

is another basis of  $\mathcal{A}_{\mathbf{C}_N^{II}}$ . Both bases are related by

$$\begin{aligned} \mathbf{B}_0^I &= \mathbf{B}_0^{II} = \mathbf{I}, \quad \mathbf{B}_1^{II} = 2\mathbf{B}_1^I, \\ \mathbf{B}_j^{II} &= \mathbf{B}_{j-2}^{II} + 2\mathbf{B}_j^I \quad (j = 2, \dots, N-1), \end{aligned} \quad (6.2)$$

where the last equation follows by  $U_j = U_{j-2} + 2T_j$ . Now we have by (4.7) that

$$\mathbf{M}_N = \sum_{k=0}^{N-1} \alpha_k \mathbf{B}_k^{II}$$

with

$$\boldsymbol{\alpha} = \mathbf{G}^{-1} \boldsymbol{\beta}^{II}, \quad \mathbf{G} := \left( \langle \mathbf{B}_j^{II}, \mathbf{B}_k^{II} \rangle \right)_{j,k=0}^{N-1}, \quad \boldsymbol{\beta}^{II} := \left( \langle \mathbf{T}'_N \mathbf{T}_N, \mathbf{B}_j^{II} \rangle \right)_{j=0}^{N-1}.$$

Clearly, we are not interested in  $\mathbf{M}_N$  itself, but in the diagonal matrix  $\operatorname{diag} \mathbf{d}$  with

$$\mathbf{M}_N = (\mathbf{C}_N^{II})' (\operatorname{diag} \mathbf{d}) \mathbf{C}_N^{II}. \quad (6.3)$$

If  $\boldsymbol{\alpha}$  is known, then we obtain  $\operatorname{diag} \mathbf{d}$  by

$$\operatorname{diag} \mathbf{d} = \sum_{k=0}^{N-1} \alpha_k \mathbf{C}_N^{II} \mathbf{B}_k^{II} (\mathbf{C}_N^{II})' = \sum_{k=0}^{N-1} \alpha_k \operatorname{diag}(U_k(c_l))_{l=0}^{N-1},$$

i.e. by definition of  $U_k$  by

$$d_0 = \sum_{k=0}^{N-1} (k+1) \alpha_k, \quad d_k = \frac{\hat{\alpha}_k}{\sin \frac{k\pi}{N}} \quad (k = 1, \dots, N-1), \quad (6.4)$$

$$(\hat{\alpha}_k)_{k=1}^{N-1} = \tilde{\mathbf{S}}_{N-1}^I (\alpha_{k-1})_{k=1}^{N-1}. \quad (6.5)$$

Thus the construction of  $\mathbf{d}$  from given  $\boldsymbol{\alpha}$  requires  $O(N \log N)$  arithmetical operations. It remains to find an efficient construction of the coefficient vector  $\boldsymbol{\alpha}$ . Therefore we use the following lemmata.

**Lemma 6.1.** For  $\mathbf{G} := \left( \langle \mathbf{B}_j^{II}, \mathbf{B}_k^{II} \rangle \right)_{j,k=0}^{N-1}$ , it holds that

$$\mathbf{G}^{-1} = \frac{1}{2N} \begin{pmatrix} 3 & 0 & -1 & & & \\ 0 & 2 & 0 & -1 & & \\ -1 & 0 & 2 & 0 & -1 & \\ \ddots & \ddots & \ddots & \ddots & \ddots & \\ & -1 & 0 & 2 & 0 & -1 \\ & & -1 & 0 & 3 & -2 \\ & & & -1 & -2 & \frac{3N-2}{N} \end{pmatrix},$$

such that the vector multiplication with  $\mathbf{G}^{-1}$  can be computed with  $O(N)$  arithmetical operations.

**Proof:** We show that  $\mathbf{G}^{-1} \mathbf{G} = \mathbf{I}$ . By definition, we have for the  $(k, l)$ -th entry of  $\mathbf{G}$  that

$$\begin{aligned} g_{k,l} &= \langle \mathbf{B}_k^{II}, \mathbf{B}_l^{II} \rangle = \text{tr} \left( (\mathbf{B}_l^{II})' \mathbf{B}_k^{II} \right) \\ &= \sum_{j=0}^{N-1} U_l(c_j) U_k(c_j). \end{aligned}$$

Then we get for  $l = 2, \dots, N-3$  that

$$2g_{k,l} - g_{k,l-2} - g_{k,l+2} = \sum_{j=0}^{N-1} U_k(c_j) (2U_l(c_j) - U_{l-2}(c_j) - U_{l+2}(c_j))$$

and further since

$$2U_l(x) - U_{l-2}(x) - U_{l+2}(x) = 4(1-x^2)U_l(x)$$

and by definition of  $U_l$  that

$$2g_{k,l} - g_{k,l-2} - g_{k,l+2} = 4 \sum_{j=1}^{N-1} \sin \frac{(k+1)j\pi}{N} \sin \frac{(l+1)j\pi}{N} = 2N \delta_{k,l}$$

$(k = 0, \dots, N-1; l = 2, \dots, N-3)$ . Straightforward computation for  $l = 0, 1, N-2, N-1$  completes the proof. ■

Set

$$\boldsymbol{\beta}^I := \left( \langle \mathbf{T}'_N \mathbf{T}_N, \mathbf{B}_j^I \rangle \right)_{j=0}^{N-1}.$$

Then we obtain by the recurrence relation (6.2) that

$$\beta_0^{II} = \beta_0^I, \quad \beta_1^{II} = 2\beta_1^I, \quad \beta_k^{II} = \beta_{k-2}^{II} + \beta_k^I \quad (k = 2, \dots, N-1). \quad (6.6)$$

Thus we can compute  $\beta^{II}$  from  $\beta^I$  with  $O(N)$  additions. The following construction of  $\beta^I$  is based on an idea of E.E. Tyrtyshnikov [21]. We split  $\mathbf{T} = \mathbf{T}_N = (t_{j-k})_{j,k=0}^{N-1}$  into a lower and an upper triangular Toeplitz matrix

$$\mathbf{T} = \mathbf{T}_L + \mathbf{T}_R$$

with diagonal entries  $t_0/2$ . Then we obtain that

$$\beta_k^I = \langle \mathbf{T}'_L \mathbf{T}_L, \mathbf{B}_k^I \rangle + \langle \mathbf{T}'_R \mathbf{T}_R, \mathbf{B}_k^I \rangle + \langle \mathbf{T}'_L \mathbf{T}_R + \mathbf{T}'_R \mathbf{T}_L, \mathbf{B}_k^I \rangle \quad (k = 0, \dots, N-1). \quad (6.7)$$

We consider the summands on the right-hand side. The matrix  $\mathbf{T}'_L \mathbf{T}_R + \mathbf{T}'_R \mathbf{T}_L$  is a symmetric Toeplitz matrix. The matrices  $\mathbf{T}'_L \mathbf{T}_L$  and  $\mathbf{T}'_R \mathbf{T}_R$  have lost their Toeplitz structure. For  $\mathbf{A} \in \mathbb{R}^{N,N}$ , we introduce the vectors  $\mathbf{s}(\mathbf{A}) = (s_k(\mathbf{A}))_{k=0}^{N-1}$ ,  $\mathbf{h}(\mathbf{A}) = (h_k(\mathbf{A}))_{k=0}^{N-1}$  and  $\tilde{\mathbf{h}}(\mathbf{A}) = (\tilde{h}_k(\mathbf{A}))_{k=0}^{N-1}$  by

$$\begin{aligned} s_k(\mathbf{A}) &:= \langle \mathbf{A}, \text{stoeplitz}((\varepsilon_k^N)^{-2} \mathbf{e}_k) \rangle \quad (k = 0, \dots, N-1), \\ h_k(\mathbf{A}) &:= \langle \mathbf{A}, \text{hankel}(\mathbf{e}_k, 0) \rangle \quad (k = 0, \dots, N-2), \\ \tilde{h}_k(\mathbf{A}) &:= \langle \mathbf{A}, \text{hankel}(0, \mathbf{e}_k) \rangle \quad (k = 0, \dots, N-2). \end{aligned}$$

Set  $h_{-1} := 0$  and  $\tilde{h}_{-1} := 0$ . Then it follows by (6.1) that

$$\langle \mathbf{A}, \mathbf{B}_k^I \rangle = (\varepsilon_k^N)^2 s_k(\mathbf{A}) + h_{k-1}(\mathbf{A}) + \tilde{h}_{k-1}(\mathbf{A}) \quad (k = 0, \dots, N-1). \quad (6.8)$$

**Lemma 6.2.** Let  $\mathbf{A} := \mathbf{T}'_L \mathbf{T}_R + \mathbf{T}'_R \mathbf{T}_L$  and let  $\mathbf{r} := \mathbf{T}'_R(t_0, t_1, \dots, t_{N-1})'$ . Then it holds that

$$\begin{aligned} s_k(\mathbf{A}) &= 2(N-k)r_k \quad (k = 0, \dots, N-1), \\ h_0(\mathbf{A}) &= r_0, \quad h_1(\mathbf{A}) = 2r_1, \\ h_k(\mathbf{A}) &= 2r_k + h_{k-2}(\mathbf{A}) \quad (k = 2, \dots, N-2), \\ \tilde{h}_k(\mathbf{A}) &= h_k(\mathbf{A}) \quad (k = 0, \dots, N-2). \end{aligned}$$

**Proof:** Since  $\mathbf{A} = \text{stoeplitz } \mathbf{r}$ , the assertion follows by definition of  $\mathbf{s}$ ,  $\mathbf{h}$  and  $\tilde{\mathbf{h}}$ . ■

**Lemma 6.3.** Let  $\mathbf{A} := \mathbf{T}'_R \mathbf{T}_R$  and let  $\mathbf{r} := ((\varepsilon_k^N)^2 t_{-k})_{k=0}^{N-1}$ . Then it holds that

$$\mathbf{s}(\mathbf{A}) = 2 \text{hankel}(Nr_0, (N-1)r_1, \dots, r_{N-1}), (0, \dots, 0, r_{N-1}) \mathbf{r}, \quad (6.9)$$

$$h_0(\mathbf{A}) = x_0, \quad h_1(\mathbf{A}) = x_1,$$

$$h_k(\mathbf{A}) = h_{k-2}(\mathbf{A}) + x_k \quad (k = 2, \dots, N-2), \quad (6.10)$$

$$\tilde{h}_0(\mathbf{A}) = y_0/2, \quad \tilde{h}_1(\mathbf{A}) = y_1,$$

$$\tilde{h}_k(\mathbf{A}) = \tilde{h}_{k-2}(\mathbf{A}) + y_k \quad (k = 2, \dots, N-2), \quad (6.11)$$

where

$$\mathbf{x} := \mathbf{T}'_R \mathbf{r},$$

$$\mathbf{y} := \text{hankel}(2(r_0, r_1, \dots, r_{N-1}), (-r_2, -r_3, \dots, -r_{N-1}, 0, 2r_{N-1})) \mathbf{r}.$$

**Proof:** Since  $\mathbf{T}_R$  is an upper triangular Toeplitz matrix, it holds that

$$a_{j,k} = a_{j-1,k-1} + r_j r_k \quad (j, k = 1 \dots, N-1). \quad (6.12)$$

Consequently, we obtain that

$$s_k = 2 \sum_{j=0}^{N-k-1} a_{j+k,j} = \sum_{j=0}^{N-1-k} (N-j-k) r_{j+k} r_j$$

which yields (6.9). The recursions (6.10) and (6.11) follow by straightforward calculation from (6.12). ■

**Lemma 6.4.** Let  $\mathbf{B} := \mathbf{T}'_R \mathbf{T}_R$  and let  $\mathbf{J} :=$  shankel  $e_{N-1}$  denote the  $N$ -th counteridentity. Then it holds that

$$\begin{aligned} s_k(\mathbf{B}) &= s_k(\mathbf{J} \mathbf{B} \mathbf{J}) \quad (k = 0, \dots, N-1), \\ h_k(\mathbf{B}) &= \tilde{h}_k(\mathbf{J} \mathbf{B} \mathbf{J}) \quad (k = 0, \dots, N-2), \\ \tilde{h}_k(\mathbf{B}) &= h_k(\mathbf{J} \mathbf{B} \mathbf{J}) \quad (k = 0, \dots, N-2), \end{aligned}$$

such that  $s(\mathbf{B}), h(\mathbf{B})$  and  $\tilde{h}(\mathbf{B})$  can be computed by (6.9) – (6.11).

**Proof:** The relations for  $s(\mathbf{B}), h(\mathbf{B})$  and  $\tilde{h}(\mathbf{B})$  follow by definition of  $\mathbf{J}$ . By

$$\mathbf{J} \mathbf{T}'_L \mathbf{T}_L \mathbf{J} = (\mathbf{J} \mathbf{T}_L \mathbf{J})' (\mathbf{J} \mathbf{T}_L \mathbf{J})$$

and since

$$\mathbf{J} \mathbf{T}_L \mathbf{J} = \text{toeplitz}((t_0, \dots, t_{N-1}), (t_0, 0, \dots, 0))$$

is an upper triangular Toeplitz matrix, we can calculate  $s(\mathbf{B}), h(\mathbf{B})$  and  $\tilde{h}(\mathbf{B})$  by (6.9) – (6.11) with  $\mathbf{A} = \mathbf{J} \mathbf{B} \mathbf{J}$  and  $\mathbf{r} = ((\varepsilon_k^N)^2 t_k)_{k=0}^{N-1}$ . ■

**Theorem 6.5.** Let  $\mathbf{T}_N := (t_{j-k})_{j,k=0}^{N-1}$ . Then the optimal preconditioner  $\mathbf{M}_N \in \mathcal{A}_{C_N^H}$  of  $\mathbf{T}'_N \mathbf{T}_N$  can be constructed with  $O(N \log N)$  arithmetical operations.

**Proof:** We compute  $\beta^I$  by (6.7), (6.8) and by the Lemmata 6.2 – 6.4. Taking into account that the multiplication of a vector with a Toeplitz matrix or a Hankel matrix requires  $O(N \log N)$  operations, the whole construction of  $\beta^I$  takes  $O(N \log N)$  arithmetical operations. From  $\beta^I$  we compute  $\beta^{II}$  by (6.6) with  $O(N)$  additions. Using Lemma 6.1, we get  $\alpha := \mathbf{G}^{-1} \beta^{II}$  at the cost of  $O(N)$  arithmetical operations. Finally, the DST-I in (6.4) to obtain  $\mathbf{d}$  from  $\alpha$  needs  $O(N \log N)$  arithmetical operations and we are done with an arithmetical complexity of  $O(N \log N)$ . ■

**Remark.** We can use similar ideas for the construction of the optimal trigonometric preconditioners with respect to the  $\mathbf{C}_N^{IV}$ ,  $\mathbf{S}_{N-1}^I$ ,  $\mathbf{S}_N^{II}$  and  $\mathbf{S}_N^{IV}$ . For the corresponding bases of  $\mathcal{A}_{O_N}$

$$\begin{aligned}\mathbf{B}_k^I &:= \mathbf{O}'_{\dim} \operatorname{diag}(T_k(c_l))_{l=0}^{\dim-1} \mathbf{O}_{\dim} \\ &= \text{stoepitz } \mathbf{e}'_k + \text{hankel } (\mathbf{u}'_k, \mathbf{v}'_k), \\ \mathbf{B}_k^{II} &:= \mathbf{O}'_{\dim} \operatorname{diag}(U_k(c_l))_{l=0}^{\dim-1} \mathbf{O}_{\dim}\end{aligned}$$

it holds that

$$\mathbf{G}^{-1} = \frac{1}{K} \begin{pmatrix} 3 & 0 & -1 & & & \\ 0 & 2 & 0 & -1 & & \\ -1 & 0 & 2 & 0 & -1 & \\ & \ddots & \ddots & \ddots & \ddots & \\ & & -1 & 0 & 2 & 0 & -1 \\ & & & -1 & 0 & g_1 & g_2 \\ & & & & -1 & g_2 & g_3 \end{pmatrix},$$

where

$\mathbf{O}_{\dim}$	$c_l$	dim	$\mathbf{u}_k$	$\mathbf{v}_k$	$K$	$g_1$	$g_2$	$g_3$
$\mathbf{C}_N^{II}$	$\cos \frac{l\pi}{N}$	$N$	$\mathbf{e}_{k-1}$	$\mathbf{e}_{k-1}$	$2N$	3	-2	$\frac{3N-2}{N}$
$\mathbf{C}_N^{IV}$	$\cos \frac{(l+1)l\pi}{2N}$	$N$	$\mathbf{e}_{k-1}$	$-\mathbf{e}_{k-1}$	$2N$	1	0	1
$\mathbf{S}_{N-1}^I$	$\cos \frac{(l+1)\pi}{N}$	$N-1$	$-\mathbf{e}_{k-2}$	$-\mathbf{e}_{k-2}$	$2N+2$	2	0	3
$\mathbf{S}_N^{II}$	$\cos \frac{(l+1)\pi}{N}$	$N$	$-\mathbf{e}_{k-1}$	$-\mathbf{e}_{k-1}$	$2N$	3	2	$\frac{3N-2}{N}$
$\mathbf{S}_N^{IV}$	$\cos \frac{(2l+1)\pi}{2N}$	$N$	$-\mathbf{e}_{k-1}$	$\mathbf{e}_{k-1}$	$2N$	1	0	1

Note that the construction of the optimal preconditioner with respect to the  $\mathbf{C}_N^{IV}$  or the  $\mathbf{S}_N^{IV}$  is especially simple.

## 7 Numerical Results

Finally, we present four examples of nonsymmetric Toeplitz systems

$$\mathbf{T}_N \mathbf{x}_N = \mathbf{b}_N, \quad (7.1)$$

for which the preconditioning of the normal equation

$$\mathbf{T}'_N \mathbf{T}_N \mathbf{x}_N = \mathbf{T}'_N \mathbf{b}_N \quad (7.2)$$

by an optimal trigonometric preconditioner  $\mathbf{M}_N = \mathbf{M}_N(\mathbf{T}'_N \mathbf{T}_N)$  of  $\mathbf{T}'_N \mathbf{T}_N$  significantly accelerates the convergence of the CG-method. We refer to the CG-method applied to the normal equation (7.2) as NCG-method. The algorithms were realized for the optimal preconditioners with respect to the DCT-II, the DST-II, the DCT-IV and

the DST-IV, respectively. Clearly, we can also use the DST-I preconditioner. Note that for symmetric Toeplitz matrices, the DST-I preconditioned CG-method for the solution of (7.1) shows a similar convergence behaviour as the DST-II preconditioned CG-method.

The fast computation of the preconditioners in the initial step and the computation of the preconditioned NCG-method were implemented in Matlab and tested on a Sun SPARCstation 20. The fast trigonometric transforms appearing both in the initialization and in the NCG-steps were taken from the C-implementation based on [3] and [19] by using the cmex-programm.

As transform length we choose  $N = 2^n$ . The right-hand side  $\mathbf{b}_N$  of (7.1) is the vector consisting of  $N$  ones. The preconditioned NCG-method starts with the zero vector and stops if  $\|r^{(j)}\|_2/\|r^{(0)}\|_2 < 10^{-7}$ , where  $r^{(j)}$  denotes the residual vector after  $j$  iterations. Our test matrices are the following four Toeplitz matrices  $\mathbf{T}_N = (t_{j-k})_{j,k=0}^{N-1}$ :

i) (see [15])

$$t_n = \begin{cases} 1/\log(2-n) & n \leq 1, \\ 1/\log(2-n) + 1/(1+n) & n = 0, \\ 1/(1+n) & n \geq 1 \end{cases}$$

ii) (see [15])

$$t_n = \begin{cases} 2 & n = 0, \\ -0.7 t_{n+1} & n \leq -1, \\ 0.9 t_{n-1} & n \geq 1. \end{cases}$$

iii) Here we use the Toeplitz matrices  $\mathbf{T}_N$  arising from the generating function

$$f(x) = x^2 e^{ix}.$$

iv)

$$t_n = \begin{cases} -1.5 & n = -1, \\ 2 & n = 0, \\ 0.5 & n = 1, \\ 0 & \text{else.} \end{cases}$$

As expected, also for large transform lengths  $N$ , the initialisation and each NCG-step can be computed very fast which reflects the arithmetic complexity of  $O(N \log N)$  for these computations. The four last columns of the following tables show the number of iterations required by the NCG-method for the different trigonometric preconditioners. The second column contains the number of iteration steps of the NCG-method without preconditioning. The columns 3 and 4 contain the numbers of iterations required by the CG-method applied to the equation

$$(\mathbf{M}_N^{-1} \mathbf{T}_N)' (\mathbf{M}_N^{-1} \mathbf{T}_N) \mathbf{x}_N = (\mathbf{M}_N^{-1} \mathbf{T}_N) \mathbf{M}_N^{-1} \mathbf{b}_N,$$

where  $\mathbf{M}_N$  denotes the optimal preconditioner of  $\mathbf{T}_N$  with respect to the DCT-II and the DST-II, respectively.

$n$	$\mathbf{I}_N$	DCT-II	DST-II	$\mathbf{C}_N^{II}$	$\mathbf{S}_N^{II}$	$\mathbf{C}_N^{IV}$	$\mathbf{S}_N^{IV}$
7	24	11	10	8	15	14	14
8	32	12	11	8	17	15	15
9	43	14	13	8	19	17	16
10	57	18	16	9	20	19	16
11	86	20	19	9	20	20	17
12	121	25	22	9	22	22	19
13	176	30	28	9	22	22	19

**Table 1:** Number of iterations for example i)

$n$	$\mathbf{I}_N$	DCT-II	DST-II	$\mathbf{C}_N^{II}$	$\mathbf{S}_N^{II}$	$\mathbf{C}_N^{IV}$	$\mathbf{S}_N^{IV}$
7	34	44	44	9	12	9	12
8	43	47	50	8	11	8	11
9	53	50	52	7	10	8	11
10	59	50	53	7	9	7	10
11	50	50	53	6	9	7	9
12	58	49	53	6	8	7	9
13	56	48	54	6	8	7	9

**Table 2:** Number of iterations for example ii)

$n$	$\mathbf{I}_N$	DCT-II	DST-I	$\mathbf{C}_N^{II}$	$\mathbf{S}_N^{II}$	$\mathbf{C}_N^{IV}$	$\mathbf{S}_N^{IV}$
5	84	72	68	29	21	47	37
6	311	124	176	52	26	84	74
7	1226	264	412	116	33	173	140
8	5220	626	980	256	40	405	302
9	**	1741	3341	664	74	1031	846

**Table 3:** Number of iterations for example iii)

$n$	$\mathbf{I}_N$	DCT-II	DST-II	$\mathbf{C}_N^{II}$	$\mathbf{S}_N^{II}$	$\mathbf{C}_N^{IV}$	$\mathbf{S}_N^{IV}$
6	88	21	37	21	9	25	24
7	201	31	67	27	8	31	30
8	435	45	125	36	8	39	38
9	929	66	294	47	9	72	65

**Table 4:** Number of iterations for example iv)

Although not all matrices in our examples fulfil the assumptions of Theorem 5.5, the preconditioning with an optimal trigonometric preconditioner of  $\mathbf{T}'_N \mathbf{T}_N$  accelerates the convergence of the NCG-method significantly.

For all examples, the preconditioning with respect to the DCT-IV and the DST-IV leads to similar numbers of iteration steps. It is easy to check that for symmetric

Toeplitz matrices  $\mathbf{T}_N$  the optimal preconditioners  $\mathbf{M}_N(\mathbf{T}_N)$  with respect to the DCT-IV and DST-IV coincide.

Except of the second example, the number of iterations differs, if we apply the preconditioners with respect to the DCT-II and the DST-II. Heuristically, this can be explained by the different structures of  $\mathcal{A}_{C_N^{II}}$  and  $\mathcal{A}_{S_N^{II}}$  and how „good” our example matrices fit into this structure. A general criterion for the choice of the optimal trigonometric preconditioner would be interesting. In this direction, it is remarkable, that the optimal preconditioner  $\mathbf{M}_N(\mathbf{T}_N) \in \mathcal{A}_{C_N^{II}}$  of the auto-covariance matrix  $\mathbf{T}_N := (\rho^{|j-k|})_{j,k=0}^{N-1}$  is “asymptotically equivalent” to  $\mathbf{T}_N$  if  $N \rightarrow \infty$ ,  $\rho \rightarrow 1$ , while the optimal preconditioner  $\mathbf{M}_N(\mathbf{T}_N) \in \mathcal{A}_{S_N^{II}}$  of  $\mathbf{T}_N$  is “asymptotically equivalent” to  $\mathbf{T}_N$  if  $N \rightarrow \infty$ ,  $\rho \rightarrow 0$  [16].

**Acknowledgement.** The authors wish to thank M. Tasche for his valuable comments concerning the relations between trigonometric transforms and Toeplitz-plus-Hankel matrices.

## References

- [1] O. Axelsson, *Iterative Solution Methods*, Cambridge University Press, Cambridge, 1996.
- [2] E. Bozzo and C. Di Fiore, On the use of certain matrix algebras associated with discrete trigonometric transforms in matrix displacement decomposition, *SIAM J. Matrix Anal. Appl.* 16: 312–326 (1995).
- [3] G. Baszenski and M. Tasche, Fast polynomial multiplication and convolution related to the discrete cosine transform, *Linear Alg. Appl.* (in print).
- [4] E. Boman and I. Koltracht, Fast transform based preconditioners for Toeplitz equations, *SIAM J. Matrix Anal. Appl.* 16: 628–645 (1995).
- [5] T.F. Chan, An optimal circulant preconditioner for Toeplitz systems, *SIAM J. Sci. Statist. Comput.* 9: 766–771 (1988).
- [6] R.H. Chang and G. Strang, Toeplitz systems by conjugate gradients with circulant preconditioner, *SIAM J. Sci. Statist. Comput.* 10: 104–119 (1989).
- [7] R.H. Chan and M.K. Ng, Conjugate gradient methods of Toeplitz systems, *SIAM Review* 38/3 (1996).
- [8] R.H. Chan and M. Yeung, Circulant preconditioners for complex Toeplitz systems, *SIAM J. Numer. Anal.* 30: 1193–1207 (1993).
- [9] R.H. Chan and M.K. Ng, Sine transform based preconditioners for symmetric Toeplitz systems, *Linear Alg. Appl.* 232: 237–259 (1996).

- [10] R.H. Chan, T.F. Chan and C.K. Wong, Cosine transform based preconditioners for total variation minimization problems in image processing, Technical Report 95-8 (64), Department of Mathematics, The Chinese University of Hong Kong, 1995.
- [11] R.H. Chan, J.G Nagy and R.J Plemmons, Circulant preconditioned Toeplitz least square problems, *SIAM J. Matrix Anal. Appl.* 15: 80–97 (1992).
- [12] R.H. Chan, J.G Nagy and R.J Plemmons, Displacement Preconditioners for Toeplitz Least Squares Iterations, *Elec. Trans. Numer. Anal.* 2: 44–56 (1994).
- [13] G. Heinig, A. Bojanczyk, Transformation techniques for Toeplitz and Toeplitz-plus-Hankel matrices, I. Transformations, *Linear Alg. Appl.* (submitted).
- [14] G. Heinig and K. Rost, Representations of Toeplitz-plus-Hankel matrices using trigonometric transforms with application to fast matrix-vector multiplication, preprint, December 1996.
- [15] T. Ku and C. Kuo, Spectral Properties of Preconditioned Rational Toeplitz Matrices: the Nonsymmetric Case, *SIAM J. Matrix Anal. Appl.* 14: 521–544 (1993).
- [16] K.R. Rao and P. Yip, *Discrete Cosine Transform*, Academic Press, Boston, 1990.
- [17] Y. Saad, Computational Fluid Dynamics, Lecture Series, van Karman Institute for Fluid Dynamics, 1994.
- [18] V. Strela and E.E. Tyrtyshnikov, Some generalisations of circulant preconditioners, Matrix Methods and Algorithms, IVM RAN Moscow: 66–73, 1990.
- [19] G. Steidl and M. Tasche, A polynomial approach to fast algorithms for discrete Fourier–cosine and Fourier–sine transforms, *Math. Comp.* 56: 281–296 (1991).
- [20] S.V. Parter, On the distribution of singular values of Toeplitz matrices, *Lin. Alg. Appl.* 80: 115–130 (1986).
- [21] E.E. Tyrtyshnikov, Optimal and superoptimal circulant preconditioners, *SIAM J. Matrix Anal. Appl.* 13: 459–473 (1992).
- [22] E.E. Tyrtyshnikov, A unifying approach to some old and new theorems on distribution and clustering, *Linear Alg. Appl.* 232: 1–43 (1996).
- [23] E.E. Tyrtyshnikov, A.Y. Yeremin and N.L. Zamarashkin, Clusters - Preconditioners - Convergence, Preprint.
- [24] Z. Wang, Fast algorithms for the discrete W transform and for the discrete Fourier transform, *IEEE Trans. Acoust. Speech Signal Process.* 32: 803–816 (1984).
- [25] J. Wilkinson, *The Algebraic Eigenvalue Problem*, Clarendon Press, Oxford, 1965.