

Consiglio Nazionale delle Ricerche

**A divide and conquer algorithm for Toeplitz-like
systems: stability analysis for
the symmetric case**

P.Favati, G.Lotti, O.Menchi

IIT TR-15/2014

Technical report

Novembre 2014



Istituto di Informatica e Telematica

A divide and conquer algorithm for Toeplitz-like systems: stability analysis for the symmetric case

P. Favati G. Lotti O. Menchi

Abstract

A superfast algorithm for the solution of Toeplitz-like systems has been introduced and tested in [4]. In this paper a theoretical error analysis of the algorithm is performed for the symmetric case. The analysis produces an upper bound to the norm of the residual vector, allowing the detection of some parameters which rule the stability behavior of the algorithm. These parameters take into account both the conditioning properties of the coefficient matrices at the different levels of recursion and the magnitude of some involved matrices measured through their generators. The experimentation confirms the theoretical results, pointing out that, in general, the upper bound to the norm of the residual vector is too pessimistic.

1 Introduction

In this paper the problem of solving the linear system

$$A\mathbf{x} = \mathbf{b}, \tag{1}$$

where A is an $N \times N$ nonsingular symmetric Toeplitz-like matrix and \mathbf{b} is a vector of size N , is addressed. The class of Toeplitz-like matrices is based on the concept of displacement rank introduced in [11] and studied by many authors (see for example [7, 9, 10]). The displacement operator allows a compact representation of the matrices of this class by means of a set of generators. For the Toeplitz-like matrices, fast and superfast algorithms have been devised (see the extensive bibliography in [10]), but the question of their stability is still matter of discussion. In [4] a new superfast algorithm, called `solvegen`, has been proposed. It is based on a divide and conquer strategy which combines the solutions of two half size Toeplitz-like systems with enlarged right-hand side increased by a constant number (proportional to the displacement rank of the given matrix) of columns at each recursion step. `solvegen` avoids the explicit inversion the matrix, which is often cause of instability. In fact, the experimentation shows that it is more stable than the superfast algorithm of [2, 12] which has $O(N \log^2 N)$ complexity. This stability improvement is obtained at the expense of an increase of the complexity to $O(N \log^3 N)$ floating point operations.

We are interested in carrying out a theoretical error analysis of `solvegen` with the aim of characterizing the parameters which rule its stability behavior. Considering the difficulty of performing such an analysis in the general case, we study the stability properties of `solvegen` rewritten for the symmetric case.

Some preliminaries on Toeplitz-like matrices are included in Section 2. In Section 3, after a first version of the proposed algorithm which does not yet make use of the displacement structure, the explicit expressions for the generators of the involved matrices and the version of the algorithm which fully exploits the displacement properties are given. The stability analysis of `solvegen` is discussed in Section 4. Finally, in Section 5 the results of the numerical experiments are shown, confirming the relevance of the detected parameters in the numerical stability of the algorithm.

2 Toeplitz-like matrices

The definition of Toeplitz-like structure is based on the concept of displacement rank, which measures how close a matrix is to a Toeplitz matrix. Given an $n \times n$ matrix A , we consider the *displacement operator*

$$\nabla(A) = A - ZAZ^T, \quad (2)$$

where Z is the $n \times n$ *down-shift* matrix

$$Z = \begin{bmatrix} 0 & & & & \\ 1 & 0 & & & \\ & \ddots & \ddots & & \\ & & & 1 & 0 \end{bmatrix}.$$

The matrix A is said to be *Toeplitz-like* if the quantity $r_{\text{disp}}(A) = \text{rank} \nabla(A)$ (called *displacement rank*) is small with respect to n (more formally $r_{\text{disp}}(A) = O(1)$ for $n \rightarrow \infty$). The set of Toeplitz-like matrices, unlike the set of Toeplitz matrices, is closed under the operations of multiplication and inversion. Let $r_{\text{disp}}(A) = \rho$, then

$$\nabla(A) = C\Sigma D^T, \quad (3)$$

for suitable $n \times \rho$ matrices C and D and a signature $\rho \times \rho$ matrix Σ , i.e. a diagonal matrix with diagonal elements σ_i either 1 or -1 . The matrices C , Σ and D are called *generators* of A . Denoting by \mathbf{c}_i and \mathbf{d}_i the columns of C and D respectively, then

$$\nabla(A) = \sum_{i=1}^{\rho} \sigma_i \mathbf{c}_i \mathbf{d}_i^T. \quad (4)$$

For example, the representation of a Toeplitz matrix A of elements $a_{i,j}$ with $a_{1,1} \neq 0$ is

$$\nabla(A) = \mathbf{c}_1 \mathbf{e}_1^T + \mathbf{e}_1 \mathbf{d}_2^T, \quad \text{with} \quad \mathbf{c}_1 = A \mathbf{e}_1, \quad \mathbf{d}_2 = A^T \mathbf{e}_1 - a_{11} \mathbf{e}_1, \quad (5)$$

i.e. \mathbf{c}_1 is the first column of A and \mathbf{d}_2^T is the first row of A without the first component. Hence, the displacement rank of a Toeplitz matrix is $\rho = 2$, except in some special case where $\rho = 1$.

From (2) and (3) it follows that

$$\nabla(A^T) = \nabla(A)^T = D\Sigma C^T, \quad (6)$$

hence the generators of A^T are obtained by swapping the generators of A . In the case of a symmetric matrix A , we have $D = C$ and

$$\nabla(A) = C\Sigma C^T. \quad (7)$$

For example, if the Toeplitz matrix whose displacement is given in (5) is symmetric and $\alpha = a_{1,1} \neq 0$, denoting $\beta = \text{sgn}(\alpha)/\sqrt{|\alpha|}$, formula (7) holds with

$$C = \beta [\mathbf{c}_1, \mathbf{c}_2], \quad \Sigma = \text{diag}(\text{sgn}(\alpha), -\text{sgn}(\alpha)), \quad \text{where } \mathbf{c}_2 = \mathbf{c}_1 - \alpha \mathbf{e}_1.$$

The generators enable us to represent a Toeplitz-like matrix as the sum of products of lower and upper triangular Toeplitz factors. Denoting by $L(\mathbf{s})$ the lower triangular Toeplitz matrix whose first column is \mathbf{s} and by $L^T(\mathbf{s})$ the upper triangular Toeplitz matrix whose first row is \mathbf{s} , then

$$A = \sum_{i=1}^{\rho} \sigma_i L(\mathbf{c}_i) L^T(\mathbf{d}_i), \quad (8)$$

with $\mathbf{d}_i = \mathbf{c}_i$ in the symmetric case. In this paper we assume that the coefficient matrix A of (1) is not explicitly given but only represented through its generators.

Formula (8) can be exploited to compute the product of a Toeplitz-like matrix A by a vector \mathbf{v} . In fact

$$A\mathbf{v} = \sum_{i=1}^{\rho} \sigma_i L(\mathbf{c}_i) L^T(\mathbf{d}_i) \mathbf{v},$$

and the product is obtained by multiplying first upper and then lower triangular Toeplitz matrices by vectors. If n is large, it is worthwhile to compute these products embedding the triangular Toeplitz matrices into circulant matrices and using FFT. In this way the product $A\mathbf{v}$ has a computational cost of $O(n \log n)$ for $n \rightarrow \infty$ (see [5] for the sketch of a function `prod` that can be used to multiply Toeplitz-like matrices by vectors).

When a new method for solving problem (1) is proposed, it is a good practice to look for conditions that ensure the stability of the method. If the elements of A are given through the generators, testing these conditions may require the explicit construction of some or all the elements of A by multiplying A with vectors of the canonical basis. For example, it is easy to verify if A is symmetric,

but if we want to test the diagonal dominance we have to compute the elements of A which in the symmetric case are given by

$$a_{k,j} = a_{j,k} = \sum_{i=1}^{\rho} \sigma_i \sum_{h=1}^j c_{h,i} c_{k-j+h,i}, \quad \text{for } j \leq k.$$

The decomposition (3) of $\nabla(A)$, and consequently the representation of A by means of the generators, is not unique. An important representation is the *orthogonal* one, obtained by computing the SVD decomposition $\nabla(A) = UWV^T$, where W is the $\rho \times \rho$ diagonal matrix of the nonzero singular values $w_1 \geq \dots \geq w_\rho > 0$ and U and V are $n \times \rho$ matrices with orthogonal columns. If A is not symmetric, (3) is given by

$$\nabla(A) = \widehat{C} \widehat{\Sigma} \widehat{D}^T, \quad \text{where } \widehat{C} = UW^{1/2}, \quad \widehat{\Sigma} = I, \quad \widehat{D} = VW^{1/2}.$$

If A is symmetric, let $\widehat{\Sigma}$ be the signature matrix such that $V = U \widehat{\Sigma}$. Then (3) is given by

$$\nabla(A) = \widehat{C} \widehat{\Sigma} \widehat{D}^T, \quad \text{where } \widehat{C} = \widehat{D} = UW^{1/2}.$$

In [3, 6] the stability of a method for solving a linear system, whose Toeplitz-like matrix A is not explicitly given but only represented through its generators, is recognized depending on how large the generators are with respect to the magnitude of A . So, when stability is analyzed, we suggest to consider the function

$$\psi(A) = \sum_{i=1}^{\rho} \|\mathbf{c}_i\|_1 \|\mathbf{d}_i\|_1, \quad (9)$$

which verifies $\|\nabla(A)\|_1 \leq \psi(A)$. Since for any vector \mathbf{s} it is $\|L(\mathbf{s})\|_1 = \|\mathbf{s}\|_1$ then from (8) it follows that

$$\|A\|_1 \leq \psi(A). \quad (10)$$

Orthogonal generators have in general smaller 1-norms; if they are used for evaluating the function $\psi(A)$, a more accurate measure of $\|A\|_1$ is often obtained. For this reason, we will resort to the replacing of the generators with their orthogonal counterparts if the latter ones give a smaller value of $\psi(A)$. This replacement does not contribute to the global computational cost, since the cost of obtaining the orthogonal representation from a generic one is $O(n\rho^2)$ for $n \rightarrow \infty$, lower than the cost of the matrix by vector product.

Notation: A matrix is denoted by an upper-case letter, possibly followed by indices, and its columns are denoted by the corresponding lower-case and bold-face letter, followed by the index of the matrix and the index of the column. For example, $\mathbf{c}_{F,i}$ is the i th column of the matrix C_F and $\mathbf{d}_{12,i}$ is the i th column of the matrix D_{12} . Moreover, to simplify the list of the input parameters of the programs, the notation $\{A\}$ will be used to indicate the set $\{n, \rho, C, \Sigma, D\}$. For simplicity the subscript 1 is dropped from the 1-norm for both vectors and matrices.

3 The algorithm

As already said, the coefficient matrix A in (1) is assumed to be a nonsingular symmetric Toeplitz-like matrix, not given explicitly but only through its representation, which we assume to be orthogonal. In [4] a divide and conquer superfast algorithm, called `solvegen`, has been described for solving system (1) in the general Toeplitz-like (i.e. possibly nonsymmetric) case and we rewrite here the algorithm for the symmetric case.

We assume $N = 2^p n_e$, where n_e represents the size of the systems (called *elementary*) to be solved at the last level of the recursion, with the restriction that $n_e \geq r_{\text{disp}}(A) = \rho$. At each recursion level, the algorithm combines the solutions of two half size Toeplitz-like systems with enlarged right-hand side and makes use of the displacement properties of the involved matrices.

Before writing the code of `solvegen` it is advisable to sketch the code implementing the divide and conquer strategy in terms of matrices.

3.1 The divide and conquer strategy

Initially we set $A^{(0)} = A$ and $B^{(0)} = \mathbf{b}$. At the k th level of recursion, $k = 0, \dots, p$, there are 2^k symmetric systems of size $n_k = 2^{p-k} n_e$ to be solved. Let

$$A^{(k)} X^{(k)} = B^{(k)} \quad (11)$$

be one of these systems. For $k \leq p-1$, let

$$A^{(k)} = \begin{bmatrix} A_{11}^{(k)} & A_{12}^{(k)} \\ A_{21}^{(k)} & A_{22}^{(k)} \end{bmatrix}, \quad X^{(k)} = \begin{bmatrix} \overline{X}^{(k)} \\ \underline{X}^{(k)} \end{bmatrix}, \quad B^{(k)} = \begin{bmatrix} \overline{B}^{(k)} \\ \underline{B}^{(k)} \end{bmatrix}, \quad (12)$$

with $A_{11}^{(k)}$ and $A_{22}^{(k)}$ symmetric and $A_{21}^{(k)} = A_{12}^{(k)T}$. The solution of (11) can be found by computing

$$\begin{aligned} [U, F] &= A_{11}^{(k)-1} [\overline{B}^{(k)}, A_{12}^{(k)}], \\ Y &= \underline{B}^{(k)} - A_{21}^{(k)} U, \quad S = A_{22}^{(k)} - A_{21}^{(k)} F, \\ \underline{X}^{(k)} &= S^{-1} Y, \quad \overline{X}^{(k)} = U - F \underline{X}^{(k)}, \end{aligned} \quad (13)$$

i.e. by solving two subsystems, the first one with matrix $A_{11}^{(k)}$ and an increased number of right-hand sides, the second one with the Schur complement S of $A_{11}^{(k)}$ without increase of right-hand sides.

The following recursive routine `solvemat` shows how the computation can be performed. The first time the function is applied with $k = 0$ to matrix $A^{(0)}$ and right-hand side $B^{(0)}$. Subsequent calls are made to matrices of halved dimension, until the preassigned size n_e is reached. At the k th level of recursion, with $k = 0, \dots, p$, we assume that all the 2^k subsystems of the form (11) have a nonsingular coefficient matrix $A^{(k)}$. At the last level of recursion there are 2^p

elementary systems with matrices of size n_e . These systems are solved directly, by means of a routine `elem` (implementing, for example, Gauss method). The initial partition of the matrix $A^{(k)}$ and of the right-hand side $B^{(k)}$ is performed according to (12) (in the following code the colon notation is used).

```

function  $X^{(k)} = \text{solvemat}(n_k, A^{(k)}, B^{(k)})$ 
%      computes recursively the solution of the system  $A^{(k)} X = B^{(k)}$  of size  $n_k$ 
      if  $k = p$ ,  $X^{(k)} = \text{elem}(A^{(k)}, B^{(k)})$ 
      else
%          partition
           $n = n_k$ ;    $r_1 = (1 : n/2)$ ;    $r_2 = (n/2 + 1 : n)$ ;
           $A_{11} = A^{(k)}(r_1, r_1)$ ;    $A_{12} = A^{(k)}(r_1, r_2)$ ;
           $A_{21} = A^{(k)}(r_2, r_1)$ ;    $A_{22} = A^{(k)}(r_2, r_2)$ ;
           $\overline{B} = B^{(k)}(r_1, :)$ ;    $\underline{B} = B^{(k)}(r_2, :)$ ;
%          recursion
           $[U, F] = \text{solvemat}(n/2, A_{11}, [\overline{B}, A_{12}])$ ;
           $Y = \underline{B} - A_{21}U$ ;    $S = A_{22} - A_{21}F$ ;
           $\underline{X} = \text{solvemat}(n/2, S, Y)$ ;    $\overline{X} = U - F\underline{X}$ ;
           $X^{(k)} = \begin{bmatrix} \overline{X} \\ \underline{X} \end{bmatrix}$ ;
      end

```

It is easy to recognize that the divide and conquer strategy used for `solvemat` leads to a block \mathcal{LU} factorization of A . In general, an excessive growth of the elements of F can be an important factor for the instability of the method. If the matrix A has some special structural properties, then all the matrices F arising at the different levels can be nicely bounded. For example, if A is positive definite, then $A^{(k)}$ is positive definite for any k and $\|F\|_2 \leq \sqrt{\kappa_2(A^{(k)})}$, where $\kappa_2(A^{(k)})$ is the condition number of $A^{(k)}$. If $A^{(k)}$ has diagonal dominance, then $\|F\| \leq 1$.

3.2 Generators of the matrices used in the recursion

We give now the representation, in terms of generators, of the matrices involved in the computation. It is important to have at any recursion level the representation with the minimal number of generators. When writing the generators of one of the matrices required at recursion level k , for simplicity of notation the level superscript k is omitted, so

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix},$$

with A_{11} and A_{22} symmetric and $A_{21} = A_{12}^T$. Let C and Σ satisfy (7) and partition C compatibly

$$C = \begin{bmatrix} \bar{C} \\ \underline{C} \end{bmatrix}.$$

For simplicity, we still denote by Z be the down-shift matrix of order $n/2$ and \mathbf{e}_i the i th canonical vector of length $n/2$.

The following scalar and vectors are computed by means of (8)

$$\begin{aligned} \alpha &= \mathbf{e}_{n/2}^T A_{11} \mathbf{e}_{n/2} = \sum_{i=1}^{\rho} \sigma_i \bar{\mathbf{c}}_i^T \bar{\mathbf{c}}_i, & \mathbf{v}_1 &= Z A_{11} \mathbf{e}_{n/2} = Z \sum_{i=1}^{\rho} \sigma_i L(\bar{\mathbf{c}}_i) J \bar{\mathbf{c}}_i, \\ \mathbf{v}_2 &= Z A_{21} \mathbf{e}_{n/2} = Z \sum_{i=1}^{\rho} \sigma_i L(\underline{\mathbf{c}}_i) J \bar{\mathbf{c}}_i, \end{aligned} \quad (14)$$

where J is the reversal matrix.

- The displacements of the blocks are

$$\begin{aligned} \nabla(A_{11}) &= \bar{C} \Sigma \bar{C}^T, \\ \nabla(A_{12}) &= \bar{C} \Sigma \underline{C}^T + \mathbf{v}_1 \mathbf{e}_1^T, \\ \nabla(A_{21}) &= \nabla(A_{12})^T = \underline{C} \Sigma \bar{C}^T + \mathbf{e}_1 \mathbf{v}_1^T, \\ \nabla(A_{22}) &= \underline{C} \Sigma \underline{C}^T + \mathbf{v}_3 \mathbf{e}_1^T + \mathbf{e}_1 \mathbf{v}_2^T, \quad \mathbf{v}_3 = \alpha \mathbf{e}_1 + \mathbf{v}_2. \end{aligned} \quad (15)$$

Then

$$\nabla(A_{21}) = C_{21} \Sigma_{21} D_{21}^T, \quad \text{where } C_{21} = [\underline{C}, \mathbf{e}_1], \quad D_{21} = [\bar{C}, \mathbf{v}_1], \quad \Sigma_{21} = \begin{bmatrix} \Sigma & \\ & 1 \end{bmatrix}.$$

In general

$$r_{\text{disp}}(A_{11}) = \rho, \quad r_{\text{disp}}(A_{12}) = r_{\text{disp}}(A_{21}) = \rho + 1, \quad r_{\text{disp}}(A_{22}) = \rho + 2.$$

- The matrix $F = A_{11}^{-1} A_{12}$ is not symmetric. Its displacement is

$$\nabla(F) = F - Z F Z^T = A_{11}^{-1} (A_{12} - A_{11} Z F Z^T).$$

Since

$$A_{12} - A_{11} Z F Z^T = \nabla(A_{12}) + Z A_{12} Z^T - (\nabla(A_{11}) + Z A_{11} Z^T) Z F Z^T,$$

using $Z^T Z = I - \mathbf{e}_{n/2} \mathbf{e}_{n/2}^T$, we have

$$\begin{aligned} Z A_{12} Z^T - Z A_{11} Z^T Z F Z^T &= Z A_{12} Z^T - Z A_{11} F Z^T + Z A_{11} \mathbf{e}_{n/2} \mathbf{e}_{n/2}^T F Z^T \\ &= \mathbf{v}_1 \mathbf{e}_{n/2}^T F Z^T, \end{aligned}$$

and

$$\nabla(F) = A_{11}^{-1}(\underline{C} \Sigma \underline{C}^T + \mathbf{v}_1 \mathbf{e}_1^T - \underline{C} \Sigma \underline{C}^T Z F Z^T + \mathbf{v}_1 \mathbf{e}_{n/2}^T F Z^T).$$

Then

$$\nabla(F) = C_+ \Sigma D_+^T + \mathbf{v}_+ \mathbf{f}^T,$$

where

$$C_+ = A_{11}^{-1} \underline{C}, \quad D_+ = \underline{C} - Z A_{21} A_{11}^{-1} Z^T \underline{C}, \quad \mathbf{v}_+ = A_{11}^{-1} \mathbf{v}_1, \quad \mathbf{f} = \mathbf{e}_1 + Z A_{21} A_{11}^{-1} \mathbf{e}_{n/2}.$$

Setting $K = [Z^T \underline{C}, -\mathbf{e}_{n/2}]$ and $G = A_{11}^{-1} K$, we have

$$\nabla(F) = C_F \Sigma_F D_F, \quad (16)$$

where

$$C_F = A_{11}^{-1} D_{21}, \quad D_F = C_{21} - Z A_{21} G, \quad \Sigma_F = \Sigma_{21}.$$

In general $r_{\text{disp}}(F) = \rho + 1$.

- The Schur complement of A_{11} is $S = A_{22} - A_{21} F$. It is symmetric and its displacement is

$$\nabla(S) = \nabla(A_{22}) - \nabla(A_{21} F).$$

It is

$$\nabla(A_{22}) = \underline{C} \Sigma \underline{C}^T + \mathbf{v}_3 \mathbf{e}_1^T + \mathbf{e}_1 \mathbf{v}_2^T = \underline{C} \Sigma D_+^T + H_1, \quad (17)$$

where

$$H_1 = \underline{C} \Sigma \underline{C}^T Z F Z^T + \mathbf{v}_3 \mathbf{e}_1^T + \mathbf{e}_1 \mathbf{v}_2^T,$$

and

$$\nabla(A_{21} F) = A_{21} F - Z A_{21} F Z^T = A_{21} \nabla(F) + (A_{21} Z - Z A_{21}) F Z^T. \quad (18)$$

Since

$$A_{21} Z - Z A_{21} = A_{21} Z - Z A_{21} (Z^T Z + \mathbf{e}_{n/2} \mathbf{e}_{n/2}^T) = \nabla(A_{21}) Z - \mathbf{v}_2 \mathbf{e}_{n/2}^T, \quad (19)$$

we have

$$\nabla(A_{21} F) = A_{21} C_+ \Sigma D_+^T + H_2, \quad (20)$$

where

$$H_2 = \nabla(A_{21}) Z F Z^T + A_{21} \mathbf{v}_+ \mathbf{f}^T - \mathbf{v}_2 \mathbf{e}_{n/2}^T F Z^T.$$

Now

$$\begin{aligned} H_1 - H_2 &= \mathbf{v}_3 \mathbf{e}_1^T + \mathbf{e}_1 \mathbf{v}_2^T - \mathbf{e}_1 \mathbf{v}_1^T Z F Z^T - A_{21} \mathbf{v}_+ \mathbf{f}^T + \mathbf{v}_2 \mathbf{e}_{n/2}^T F Z^T \\ &= \mathbf{g} \mathbf{f}^T, \quad \text{where } \mathbf{g} = \mathbf{v}_3 - A_{21} \mathbf{v}_+. \end{aligned}$$

From (17) and (20) it follows that

$$\nabla(S) = (\underline{C} - A_{21} C_+) \Sigma D_+^T + \mathbf{g} \mathbf{f}^T = C_S \Sigma_S D_S^T,$$

where

$$C_S = [\underline{C}, \mathbf{v}_3] - A_{21}C_F, \quad \Sigma_S = \Sigma_F, \quad D_S = D_F. \quad (21)$$

Actually, the matrix $\underline{C} - A_{21}C_+$ has ρ columns but rank $\rho - 1$. In fact, if we multiply it by the vector $\Sigma \overline{C}^T \mathbf{e}_1$ we obtain the vector

$$(\underline{C} - A_{21}C_+) \Sigma \overline{C}^T \mathbf{e}_1 = \nabla(A_{21})\mathbf{e}_1 - \mathbf{e}_1 \mathbf{v}_1^T \mathbf{e}_1 - A_{21}A_{11}^{-1} \nabla(A_{11})\mathbf{e}_1,$$

which is null because $\mathbf{v}_1^T \mathbf{e}_1 = 0$, $\nabla(A_{21})\mathbf{e}_1 = A_{21}\mathbf{e}_1$ and $\nabla(A_{11})\mathbf{e}_1 = A_{11}\mathbf{e}_1$. Then $r_{\text{disp}}(S) = \rho$ and the representation (21) is not minimal. Moreover, we expect a symmetric representation for S , while the generators C_S and D_S given in (21) do not coincide. For this reason the decomposition (16) is replaced by its orthogonal decomposition $\nabla(S) = \widehat{C}_S \widehat{\Sigma}_S \widehat{D}_S$.

For what concerns the function ψ , since

$$|\alpha| \leq \sum_{i=1}^{\rho} \|\bar{\mathbf{c}}_i\|^2, \quad \|\mathbf{v}_1\| \leq \sum_{i=1}^{\rho} \|\bar{\mathbf{c}}_i\|^2, \quad \|\mathbf{v}_2\| \leq \sum_{i=1}^{\rho} \|\mathbf{c}_i\| \|\bar{\mathbf{c}}_i\|,$$

we have

$$\begin{aligned} \psi(A_{11}) &\leq \sum_{i=1}^{\rho} \|\bar{\mathbf{c}}_i\|^2 \leq \sum_{i=1}^{\rho} \|\mathbf{c}_i\|^2 = \psi(A), \\ \psi(A_{21}) = \psi(A_{12}) &\leq \sum_{i=1}^{\rho} (\|\bar{\mathbf{c}}_i\| \|\mathbf{c}_i\| + \|\bar{\mathbf{c}}_i\|^2) \leq \sum_{i=1}^{\rho} \|\bar{\mathbf{c}}_i\| \|\mathbf{c}_i\| \leq \sum_{i=1}^{\rho} \|\mathbf{c}_i\|^2 = \psi(A), \\ \psi(A_{22}) &\leq \sum_{i=1}^{\rho} (\|\mathbf{c}_i\|^2 + 2\|\bar{\mathbf{c}}_i\| \|\mathbf{c}_i\| + \|\bar{\mathbf{c}}_i\|^2) = \sum_{i=1}^{\rho} \|\mathbf{c}_i\|^2 \leq \psi(A). \end{aligned}$$

Analogously

$$\|\nabla(A_{21}) - \mathbf{v}_2 \mathbf{e}_{n/2}^T\| \leq \psi(A),$$

and from (18) and (19) it follows that

$$\nabla(S) = \nabla(A_{22}) - \nabla(A_{21}F) = \nabla(A_{22}) - A_{21} \nabla(F) - (\nabla(A_{21}) - \mathbf{v}_2 \mathbf{e}_{n/2}^T) F Z^T,$$

then

$$\psi(S) \leq \psi(A_{22}) + \|A_{21}\| \psi(F) + \psi(A) \|F\| \leq 2 \psi(A) (1 + \psi(F)). \quad (22)$$

3.3 The algorithm using the generators

We are now ready to rewrite the code of the function `solvemat` in terms of the generators. The following functions are required:

- `elemgen` to solve the elementary systems with symmetric matrices of size n_e .
- `replace` to replace the generators using the SVD.
- `prod` to multiply a Toeplitz-like matrix by vectors.

For stability reason, also the generators of A_{21} and F are replaced by the orthogonal ones, if these result in a smaller value of $\psi(A_{21})$ and $\psi(F)$, respectively.

```

function  $X^{(k)} = \text{solvegen}(\{A^{(k)}\}, B^{(k)})$ 
%   computes recursively the solution of the system  $A^{(k)} X = B^{(k)}$ 
%    $A^{(k)}$  has size  $n_k$ , displacement rank  $\rho$  and  $\nabla(A^{(k)}) = C^{(k)} \Sigma^{(k)} C^{(k)T}$ .
if  $n = n_e$ ,  $X^{(k)} = \text{elemgen}(A^{(k)}, B^{(k)})$ ;
else
     $n = n_k$ ;
%   partition into matrices of size  $n_{k+1}$ 
 $r_1 = (1 : n/2)$ ;  $r_2 = (n/2 + 1 : n)$ ;  $\overline{B} = B^{(k)}(r_1, :)$ ;  $\underline{B} = B^{(k)}(r_2, :)$ ;
 $\overline{C} = C^{(k)}(r_1, :)$ ;  $\underline{C} = C^{(k)}(r_2, :)$ ;  $\Sigma = \Sigma^{(k)}$ ;
    compute the scalar  $\alpha$  and the vectors  $\mathbf{v}_1$ ,  $\mathbf{v}_2$  and  $\mathbf{v}_3$  defined in (14) and (15),
 $C_{21} = [\underline{C}, \mathbf{e}_1]$ ;  $D_{21} = [\overline{C}, \mathbf{v}_1]$ ;  $K = [Z^T \overline{C}, -\mathbf{e}_{n/2}]$ ;  $\Sigma_{21} = \begin{bmatrix} \Sigma & 0 \\ 0 & 1 \end{bmatrix}$ ;
     $(\widehat{C}_{21}, \widehat{\Sigma}_{21}, \widehat{D}_{21}) = \text{replace}(C_{21}, \Sigma_{21}, D_{21})$ ;
%   first recursion substep, with  $\nabla(A_{11}) = \overline{C} \Sigma \overline{C}^T$ 
 $[U, C_F, G] = \text{solvegen}(\{A_{11}\}, [\overline{B}, D_{21}, K])$ ;
 $D_F = C_{21} - Z \text{prod}(\{A_{21}\}, G)$ ;
 $(\widehat{C}_F, \widehat{\Sigma}_F, \widehat{D}_F) = \text{replace}(C_F, \Sigma_{21}, D_F)$ ;
 $C_S = [\underline{C}, \mathbf{v}_3] - \text{prod}(\{A_{21}\}, C_F)$ ;
 $(\widehat{C}_S, \widehat{\Sigma}_S, \widehat{D}_S) = \text{replace}(C_S, \Sigma_S, D_S)$ ;
%   second recursion substep, with  $\nabla(S) = \widehat{C}_S \widehat{\Sigma}_S \widehat{D}_S^T$ 
 $Y = \underline{B} - \text{prod}(\{A_{21}\}, U)$ ;
 $\underline{X} = \text{solvegen}(\{S\}, Y)$ ;
 $\overline{X} = U - \text{prod}(\{F\}, \underline{X})$ ;
 $X^{(k)} = \begin{bmatrix} \overline{X} \\ \underline{X} \end{bmatrix}$ ;
end

```

The first time, the function is applied with $k = 0$, the generators of $A^{(0)} = A$ and the right-hand side $B^{(0)} = \mathbf{b}$. At the elementary level, any routine can be chosen as `elemgen`. A natural choice would be Levinson algorithm, which can be applied directly to the generators. If stability is at risk, a more stable method, like Cholesky method, can be chosen, but in this case the reconstruction of the elementary matrix from the generators is required.

In [4] the computational cost and the memory requirement of `solvegen` are given. Since the computational cost is of order $O(N \log_2^3 N)$ for $N \rightarrow \infty$, `solvegen` belongs to the class of superfast methods.

4 Analysis of stability

For the stability analysis we assume that the computations are carried out in a floating point arithmetic with unit roundoff ϵ such that $N\epsilon \ll 1$. The computed value of a variable (scalar, vector or matrix) v is denoted by \tilde{v} or by “ $fl(v)$ ”. We assume also that the quantities which appear in the bounds are not so large to invalidate a first order error analysis. For simplicity the term “ $+O(\epsilon^2)$ ”, which appears in the thesis of the theorems, is omitted in the proofs. Consequently, any expression of the form $x\tilde{y}$, where $x = O(\epsilon)$ and $\tilde{y} - y = O(\epsilon)$, is replaced by xy .

From ([8], Ch. 3) the following bound can be derived. Let A be an $n \times n$ matrix and \mathbf{u} , \mathbf{v} and \mathbf{w} be three vectors, with $\mathbf{w} = \mathbf{u} - A\mathbf{v}$. If standard multiplication techniques are used to compute $\tilde{\mathbf{w}} = fl(\mathbf{u} - fl(A\mathbf{v}))$, a matrix $\Theta_{A,\mathbf{v}}$ exists such that

$$\tilde{\mathbf{w}} = \mathbf{u} - (A + \Theta_A)\mathbf{v} + \text{diag}(\boldsymbol{\epsilon})\mathbf{w}, \quad \text{with} \quad \|\Theta_A\| \leq \epsilon n \|A\| + O(\epsilon^2), \quad (23)$$

where $\boldsymbol{\epsilon}$ is a vector whose components are bounded in modulus by ϵ . In the case A is Toeplitz-like, the product $A\mathbf{v}$ is computed by means of a function `prod` and

$$\tilde{\mathbf{w}} = fl(\mathbf{u} - \text{prod}(\{A\}, \mathbf{v})).$$

Due to the well-known stability properties of FFT [1], we assume that a bound similar to (23) holds

$$\tilde{\mathbf{w}} = \mathbf{u} - (A + \Theta_A)\mathbf{v} + \text{diag}(\boldsymbol{\epsilon})\mathbf{w}, \quad \text{with} \quad \|\Theta_A\| \leq \epsilon \gamma(n) \psi(A) + O(\epsilon^2), \quad (24)$$

where $\gamma(n) = c n \log_2 n$ (c being a function of ρ , for a proof of this bound see [5]).

The present analysis aims at detecting the parameters which play the major role in the stability of `solvegen`. Algorithm `solvegen` can be regarded as composed of two parts:

- the *outer* part, which forms the framework of the algorithm and implements the recursion (already outlined in `solvemat`) using the generators of F and S ,
- the *inner* part, which computes the generators of F and S (this part makes the essential difference of `solvegen` with respect to `solvemat`).

Following this approach, the error analysis of `solvegen` will deal first in Section 4.1 with the computation of the recursion, assuming that the errors of the computed generators are sufficiently bounded from above. Then the conditions upon which the computation of the generators meets these stability requirements will be analyzed in Section 4.2.

In this analysis we neglect the error due to the computation of the orthogonal generators by the SVD algorithm, which is commonly considered stable.

4.1 Stability of the outer part

At the k th level of recursion, with $k = 0, \dots, p$, there are 2^k systems to be solved, each one of size $n_k = 2^{p-k}n_e$. For $r = 1, \dots, 2^k$, let $\mathcal{S}_r^{(k)}$ be one of these systems. Let $A_r^{(k)}$ be its coefficient matrix, decomposed as in (12), $B_r^{(k)}$ be its right-hand side having $m_r^{(k)}$ columns and $X_r^{(k)}$ be its solution. To simplify the notation, we drop the indices k and r , denoting the system

$$AX = B, \quad \text{with} \quad A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \quad X = \begin{bmatrix} \bar{X} \\ \underline{X} \end{bmatrix}, \quad B = \begin{bmatrix} \bar{B} \\ \underline{B} \end{bmatrix}. \quad (25)$$

From $\mathcal{S}_r^{(k)}$ two subsystems $\mathcal{S}_{2r-1}^{(k+1)}$ and $\mathcal{S}_{2r}^{(k+1)}$ are derived, the first one with matrix A_{11} , the second one with the Schur matrix $S = A_{22} - A_{21}F$. Denoting by \mathbf{x}_i and \mathbf{b}_i the i th column of X and B , system (25) combines the m systems

$$A\mathbf{x}_i = \mathbf{b}_i, \quad \text{for} \quad i = 1, \dots, m.$$

To further simplify the notation, we also drop the column index i , and consider the system

$$A\mathbf{x} = \mathbf{b}, \quad \text{with} \quad \mathbf{x} = \begin{bmatrix} \bar{\mathbf{x}} \\ \underline{\mathbf{x}} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \bar{\mathbf{b}} \\ \underline{\mathbf{b}} \end{bmatrix} \quad (26)$$

to represent the i th system having $A_r^{(k)}$ as coefficient matrix, the i th column $\mathbf{b}_{r,i}^{(k)}$ of $B_r^{(k)}$ as right-hand side, and $\mathbf{x}_{r,i}^{(k)}$ as solution, for selected i, k and r .

At the last level of the recursion, for $k = p$, the size of (26) is n_e and we assume that the computed solution $\tilde{\mathbf{x}}$ satisfies

$$(A + \Delta_A)\tilde{\mathbf{x}} = \mathbf{b}, \quad \text{with} \quad \|\Delta_A\| \leq \epsilon \pi_e \psi(A) + O(\epsilon^2), \quad (27)$$

where π_e is a function of n_e , related to the method used at the elementary level (a similar bound appears in Theorem 10.4 of [8] for Cholesky method).

The framework of the outer part of `solvegen` applied to (27) is composed of the following instructions

$$\begin{aligned} \mathbf{u} &= \text{solvegen}(\{A_{11}\}, \bar{\mathbf{b}}), & \mathbf{y} &= \underline{\mathbf{b}} - \text{prod}(\{A_{21}\}, \mathbf{u}), \\ \underline{\mathbf{x}} &= \text{solvegen}(\{S\}, \mathbf{y}), & \bar{\mathbf{x}} &= \mathbf{u} - \text{prod}(\{F\}, \underline{\mathbf{x}}). \end{aligned} \quad (28)$$

Obviously, the errors of this computation depend on how well the generators \tilde{C}_F and \tilde{D}_F of F and \tilde{C}_S and $\tilde{\Sigma}_S$ of S have been determined in the inner part. Thus, in this outer part we assume that the errors of \tilde{F} and \tilde{S} are sufficiently bounded from above.

Theorem 1 *Consider the $n_k \times n_k$ system (26) for selected i, k and r . If the generators $\tilde{C}_F, \tilde{D}_F, \tilde{C}_S, \tilde{\Sigma}_S$ computed in the inner part are such that the corresponding matrices \tilde{F} and \tilde{S} satisfy*

$$\|A_{11}(\tilde{F} - F)\| + \|\tilde{S} - S'\| \leq \epsilon \omega_{k,r} \psi(A) (1 + \psi(F)) + O(\epsilon^2), \quad (29)$$

for $S' = A_{22} - A_{21}\tilde{F}$ and a suitable $\omega_{k,r}$, then the solution of system (26) computed by $\tilde{\mathbf{x}} = fl(\text{solvegen}(\{A\}, \mathbf{b}))$ satisfies

$$(A + \Delta_A)\tilde{\mathbf{x}} = \mathbf{b}, \text{ with } \|\Delta_A\| \leq \epsilon \pi_k \psi(A) + O(\epsilon^2), \quad (30)$$

for a suitable π_k not depending on the indices i and r .

Proof. The proof is by induction on k . For $k = p$, the thesis follows from (27), with $\pi_p = \pi_e$. For $k = p - 1, \dots, 0$, the quantities effectively computed by instructions (28) are

- $\tilde{\mathbf{u}} = fl(\text{solvegen}(\{A_{11}\}, \bar{\mathbf{b}}))$. By the inductive hypothesis (working to first order)

$$(A_{11} + \Delta_{A_{11}})\tilde{\mathbf{u}} = \bar{\mathbf{b}}, \text{ with } \|\Delta_{A_{11}}\| \leq \epsilon \pi_{k+1} \psi(A_{11}). \quad (31)$$

- $\tilde{\mathbf{y}} = fl(\underline{\mathbf{b}} - fl(\text{prod}(\{A_{21}\}, \tilde{\mathbf{u}})))$. From (24), setting $\gamma_{k+1} = \gamma(n_k/2)$

$$\tilde{\mathbf{y}} = \underline{\mathbf{b}} - (A_{21} + \Theta_{A_{21}})\tilde{\mathbf{u}} + \text{diag}(\epsilon)\tilde{\mathbf{y}}, \text{ with } \|\Theta_{A_{21}}\| \leq \epsilon \gamma_{k+1} \psi(A_{21}). \quad (32)$$

- $\tilde{\mathbf{x}} = fl(\text{solvegen}(\{\tilde{S}\}, \tilde{\mathbf{y}}))$. By the inductive hypothesis

$$(\tilde{S} + \Delta_S)\tilde{\mathbf{x}} = \tilde{\mathbf{y}} \text{ with } \|\Delta_S\| \leq \epsilon \pi_{k+1} \psi(S). \quad (33)$$

- $\tilde{\tilde{\mathbf{x}}} = fl(\tilde{\mathbf{u}} - fl(\text{prod}(\{\tilde{F}\}, \tilde{\tilde{\mathbf{x}}}))$. From (24)

$$\tilde{\tilde{\mathbf{x}}} = \tilde{\mathbf{u}} - (\tilde{F} + \Theta_F)\tilde{\tilde{\mathbf{x}}} + \text{diag}(\epsilon)\tilde{\tilde{\mathbf{x}}} \text{ with } \|\Theta_F\| \leq \epsilon \gamma_{k+1} \psi(F). \quad (34)$$

Replacing $\tilde{\mathbf{u}}$ from (34) into (31) we get

$$\bar{\mathbf{b}} = A_{11}(\tilde{\tilde{\mathbf{x}}} + (\tilde{F} + \Theta_F)\tilde{\tilde{\mathbf{x}}} - \text{diag}(\epsilon)\tilde{\tilde{\mathbf{x}}}) + \Delta_{A_{11}}(\tilde{\tilde{\mathbf{x}}} + \tilde{F}\tilde{\tilde{\mathbf{x}}})$$

replacing $\tilde{\mathbf{y}}$ from (33) and $\tilde{\mathbf{u}}$ from (34) into (32) we get

$$\underline{\mathbf{b}} = (\tilde{S} + \Delta_S)\tilde{\tilde{\mathbf{x}}} + A_{21}(\tilde{\tilde{\mathbf{x}}} + (\tilde{F} + \Theta_F)\tilde{\tilde{\mathbf{x}}} - \text{diag}(\epsilon)\tilde{\tilde{\mathbf{x}}}) + \Theta_{A_{21}}(\tilde{\tilde{\mathbf{x}}} + \tilde{F}\tilde{\tilde{\mathbf{x}}}) - \text{diag}(\epsilon)\tilde{S}\tilde{\tilde{\mathbf{x}}}.$$

Then

$$\begin{bmatrix} \bar{\mathbf{b}} \\ \underline{\mathbf{b}} \end{bmatrix} = \begin{bmatrix} A_{11}\tilde{\tilde{\mathbf{x}}} + A_{12}\tilde{\tilde{\mathbf{x}}} \\ A_{21}\tilde{\tilde{\mathbf{x}}} + A_{22}\tilde{\tilde{\mathbf{x}}} \end{bmatrix} + \Delta_A \begin{bmatrix} \tilde{\tilde{\mathbf{x}}} \\ \tilde{\tilde{\mathbf{x}}} \end{bmatrix} \quad (35)$$

where

$$\Delta_A = K_1 + K_2, \text{ with } K_1 = \begin{bmatrix} O & A_{11}(\tilde{F} - F) \\ O & \tilde{S} - S' + \Delta_S - \text{diag}(\epsilon)\tilde{S} \end{bmatrix},$$

$$K_2 = \begin{bmatrix} -A_{11}\text{diag}(\epsilon) + \Delta_{A_{11}} & A_{11}\Theta_F + \Delta_{A_{11}}\tilde{F} \\ -A_{21}\text{diag}(\epsilon) + \Theta_{A_{21}} & A_{21}\Theta_F + \Theta_{A_{21}}\tilde{F} \end{bmatrix}.$$

For the 1-norm of K_1 we get from (29), (33) and (22)

$$\|K_1\| \leq \epsilon (\omega_{k,r} + 2\pi_{k+1} + 1) \psi(A) (1 + \psi(F)).$$

For the 1-norm of K_2 we get from (31), (32) and (34)

$$\|K_2\| \leq \epsilon (\pi_{k+1} + 3\gamma_{k+1}) \psi(A) (1 + \psi(F)).$$

Then

$$\|\Delta_A\| \leq \epsilon \pi_k \psi(A),$$

where

$$\pi_k = (3\pi_{k+1} + 3\gamma_{k+1} + \omega_{k,r} + 1) (1 + \psi(F)). \quad (36)$$

□

4.2 Stability of the inner part

The following theorem analyzes how the generators of \tilde{F} and \tilde{S} computed by the instructions

$$\begin{aligned} C_F &= \text{solvegen}(\{A_{11}\}, D_{21}), \\ D_F &= C_{21} - Z \text{prod}(\{A_{21}\}, G), \text{ where } G = \text{solvegen}(\{A_{11}\}, K), \\ C_S &= [\underline{C}, \mathbf{v}_3] - \text{prod}(\{A_{21}\}, C_F). \end{aligned}$$

satisfy the requirement of Theorem 1.

Theorem 2 *The matrices \tilde{F} such that $\nabla(\tilde{F}) = \tilde{C}_F \Sigma_F \tilde{D}_F^T$ and \tilde{S} such that $\nabla(\tilde{S}) = \tilde{C}_S \Sigma_S \tilde{D}_S$, considered at the k th inductive step of the proof of Theorem 1, satisfy relation (29) for a suitable $\omega_{k,r}$.*

Proof. For $j = 1, \dots, \rho + 1$, consider the j th component of $\nabla(F)$, i.e. $\sigma_{F,j} \mathbf{c}_{F,j} \mathbf{d}_{F,j}^T$, with $\sigma_{F,j} = \sigma_j$ for $j = 1, \dots, \rho$ and $\sigma_{F,\rho+1} = 1$.

- The j th left generator $\mathbf{c}_{F,j}$ is computed by

$$\tilde{\mathbf{c}}_{F,j} = fl(\text{solvegen}(\{A_{11}\}, \mathbf{d}_{21})).$$

The error $\delta \mathbf{c}_{F,j} = \tilde{\mathbf{c}}_{F,j} - \mathbf{c}_{F,j}$ satisfies $\|\delta \mathbf{c}_{F,j}\| \leq \epsilon \eta_{c_j} \|\mathbf{c}_{F,j}\|$, where $\epsilon \eta_{c_j}$ is an upper bound to the relative error.

- The j th right generator $\mathbf{d}_{F,j}$ is computed by

$$\tilde{\mathbf{g}}_j = fl(\text{solvegen}(\{A_{11}\}, \mathbf{k}_j)), \quad \tilde{\mathbf{d}}_{F,j} = fl(C_{21} - Z fl(\text{prod}(\{A_{21}\}, \tilde{\mathbf{g}}_j))).$$

The error $\delta \mathbf{g}_j = \tilde{\mathbf{g}}_j - \mathbf{g}_j$ satisfies $\|\delta \mathbf{g}_j\| \leq \epsilon \eta_{g_j} \|\mathbf{g}_j\|$, where $\epsilon \eta_{g_j}$ is an upper bound to the relative error. By (24)

$$\tilde{\mathbf{d}}_{F,j} = \mathbf{c}_{21,j} - Z(A_{21} + \Theta_{A_{21}}) \tilde{\mathbf{g}}_j + \text{diag}(\epsilon_j) \mathbf{d}_{F,j}, \quad \|\Theta_{A_{21}}\| \leq \epsilon \gamma_{k+1} \psi(A_{21}).$$

Replacing $\tilde{\mathbf{g}}_i$ we get $\tilde{\mathbf{d}}_{F,j} = \mathbf{d}_{F,j} + \delta \mathbf{d}_{F,j}$, where

$$\delta \mathbf{d}_{F,j} = -Z(A_{21} \delta \mathbf{g}_j + \Theta_{A_{21}} \mathbf{g}_j) + \text{diag}(\epsilon_j) \mathbf{d}_{F,j},$$

then

$$\|\delta \mathbf{d}_{F,j}\| \leq \epsilon (\eta_{g_j} \|A_{21}\| + \gamma_{k+1} \psi(A_{21})) \|\mathbf{g}_j\| + \epsilon \|\mathbf{d}_{F,j}\|.$$

We introduce now the parameter

$$\tau_{k+1,r} = \max_{j=1, \dots, \rho+1} \frac{\|\mathbf{g}_j\|}{\|\mathbf{d}_{F,j}\|} \psi(A_{21})$$

to measure the largest cancellation error in the computation of the columns of $D_F = C_{12} - Z A_{21} G$. We get

$$\|\delta \mathbf{d}_{F,j}\| \leq \epsilon (\tau_{k+1,r} (\eta_{g_j} + \gamma_{k+1}) + 1) \|\mathbf{d}_{F,j}\|.$$

Since

$$\tilde{\mathbf{c}}_{F,j} \tilde{\mathbf{d}}_{F,j}^T - \mathbf{c}_{F,j} \mathbf{d}_{F,j}^T = \mathbf{c}_{F,j} \delta \mathbf{d}_{F,j}^T + \delta \mathbf{c}_{F,j} \mathbf{d}_{F,j}^T,$$

we get

$$\begin{aligned} \nabla(\tilde{F}) - \nabla(F) &= \sum_{j=1}^{\rho+1} \sigma_j \tilde{\mathbf{c}}_{F,j} \tilde{\mathbf{d}}_{F,j}^T - \sum_{j=1}^{\rho+1} \sigma_j \mathbf{c}_{F,j} \mathbf{d}_{F,j}^T \\ &= \sum_{j=1}^{\rho+1} \mathbf{c}_{F,j} \delta \mathbf{d}_{F,j}^T + \sum_{j=1}^{\rho+1} \delta \mathbf{c}_{F,j} \mathbf{d}_{F,j}^T. \end{aligned}$$

Setting

$$\eta_{k+1,r} = \max_{j=1, \dots, \rho+1} \max\{\eta_{c_j}, \eta_{g_j}\},$$

from (10) we have

$$\|\tilde{F} - F\| \leq \sum_{j=1}^{\rho+1} \|\mathbf{c}_{F,j}\| \|\delta \mathbf{d}_{F,j}\| + \sum_{j=1}^{\rho+1} \|\delta \mathbf{c}_{F,j}\| \|\mathbf{d}_{F,j}\| \leq \epsilon \omega_F \psi(F),$$

where

$$\omega_F = \eta_{k+1,r} (\tau_{k+1,r} + 1) + \gamma_{k+1} \tau_{k+1,r} + 1. \quad (37)$$

- Setting $Q = [\underline{C}, \mathbf{v}_3]$, from (21) we have

$$\nabla(S') = C_{S'} \Sigma_F D_{S'}^T, \quad \text{where } C_{S'} = Q - A_{21} \tilde{C}_F, \quad D_{S'} = \tilde{D}_F.$$

The generators of \tilde{S} are $\tilde{C}_S = \tilde{C}_{S'}$ and $\tilde{D}_S = \tilde{D}_{S'}$. The j th left generator $\mathbf{c}_{S,j}$ is effectively computed by

$$\tilde{\mathbf{c}}_{S,j} = fl\left(\mathbf{q}_j - fl(\text{prod}(\{A_{21}\}, \tilde{\mathbf{c}}_{F,j})\right).$$

By (24) we have

$$\tilde{\mathbf{c}}_{S,j} = \mathbf{q}_j - (A_{21} + \Theta_{A_{21}}) \tilde{\mathbf{c}}_{F,j} + \text{diag}(\epsilon_j) \mathbf{c}_{S,j}, \quad \text{with } \|\Theta_{A_{21}}\| \leq \epsilon \gamma_{k+1} \psi(A_{21}),$$

hence

$$\tilde{\mathbf{c}}_{S,j} = \mathbf{c}_{S',j} - \Theta_{A_{21}} \mathbf{c}_{F,j} + \text{diag}(\boldsymbol{\epsilon}_j) \mathbf{c}_{S,j},$$

and

$$\begin{aligned} \nabla(\tilde{S}) - \nabla(S') &= (\tilde{C}_S - C_{S'}) \Sigma_F \tilde{D}_F^T = \sum_{j=1}^{\rho+1} \sigma_j (\tilde{\mathbf{c}}_{S,j} - \mathbf{c}_{S',j}) \tilde{\mathbf{d}}_{F,j}^T \\ &= - \sum_{j=1}^{\rho+1} \sigma_j \Theta_{A_{21}} \mathbf{c}_{F,j} \mathbf{d}_{F,j}^T + \sum_{j=1}^{\rho+1} \sigma_j \text{diag}(\boldsymbol{\epsilon}_j) \mathbf{c}_{S,j} \mathbf{d}_{S,j}^T. \end{aligned} \quad (38)$$

From (10) it follows that

$$\begin{aligned} \|\tilde{S} - S'\| &\leq \sum_{j=1}^{\rho+1} \|\Theta_{A_{21}} \mathbf{c}_{F,j}\| \|\mathbf{d}_{F,j}\| + \epsilon \sum_{j=1}^{\rho+1} \|\mathbf{c}_{S,j}\| \|\mathbf{d}_{S,j}\| \\ &\leq \epsilon (\gamma_{k+1} \psi(A_{21}) \psi(F) + \psi(S)), \end{aligned}$$

and from (22) we have

$$\|\tilde{S} - S'\| \leq \epsilon \omega_S \psi(A) (1 + \psi(F)), \quad \text{where } \omega_S = \gamma_{k+1} + 2. \quad (39)$$

From (37) and (39) it follows that relation (29) holds with

$$\omega_{k,r} = \omega_F + \omega_S = (\eta_{k+1,r} + \gamma_{k+1}) (\tau_{k+1,r} + 1) + 3. \quad (40)$$

□

4.3 Stability of the whole algorithm

We are now able to investigate the stability of `solvegen`.

Theorem 3 *Let $\tilde{\mathbf{x}}$ be the solution of (1) computed by `solvegen`. Assuming that the solutions of the elementary systems satisfy (27), a matrix Δ_A exists such that*

$$(A + \Delta_A) \tilde{\mathbf{x}} = \mathbf{b}, \quad \text{with } \|\Delta_A\| \leq \epsilon \pi_0 \psi(A) + O(\epsilon^2),$$

for a suitable π_0 .

Proof. The solution of (1) computed by `solvegen` is $\tilde{\mathbf{x}} = \tilde{\mathbf{x}}_{1,1}^{(0)}$. Using the result of Theorem 1 we can give an estimate of the quantity π_0 which appears in the bound (30) for $k = 0$, i.e.

$$(A^{(0)} + \Delta_{A^{(0)}}) \tilde{\mathbf{x}}^{(0)} = \mathbf{b}^{(0)}, \quad \|\Delta_{A^{(0)}}\| \leq \epsilon \pi_0 \psi(A^{(0)}),$$

where $A^{(0)} = A$ and $\mathbf{b}^{(0)} = \mathbf{b}$. At the k th level of recursion, when the systems to be solved have $n_k = 2^{p-k}$ rows, restoring the notation with the indices of the

involved matrices and replacing $\omega_{k,r}$ in (36) with (40), the thesis of Theorem 1 becomes

$$(A_r^{(k)} + \Delta_{A_r^{(k)},i}) \tilde{\mathbf{x}}_{r,i}^{(k)} = \mathbf{b}_{r,i}^{(k)}, \quad \|\Delta_{A_r^{(k)},i}\| \leq \epsilon \pi_k \psi(A_r^{(k)}), \quad i = 1, \dots, m_r,$$

where

$$\pi_k = (1 + \psi(F_r^{(k)})) (3\pi_{k+1} + \gamma_{k+1}(\tau_{k+1,r} + 4) + \eta_{k+1,r}(\tau_{k+1,r} + 1) + 4).$$

Setting

$$\begin{aligned} \psi_k &= \max_r (1 + \psi(F_r^{(k)})), \quad \nu_k = n_k \log_2 n_k, \\ \alpha_k &= c_1 \max_r \tau_{k,r}, \quad \beta_k = c_2 \max_r \eta_{k,r} \tau_{k,r}, \end{aligned}$$

for suitable constants c_1 and c_2 , the recursive relation becomes

$$\pi_k \leq (3\pi_{k+1} + \alpha_{k+1}\nu_{k+1} + \beta_{k+1})\psi_k, \quad \text{for } k = 0, \dots, p-1.$$

Then

$$\begin{aligned} \pi_0 &\leq 3^p \pi_e \prod_{k=1}^p \psi_{k-1} + \sum_{k=1}^p (3^{k-1} (\alpha_k \nu_k + \beta_k) \prod_{i=1}^k \psi_{i-1}) \\ &\leq \Psi \left(3^p \pi_e + \sum_{k=1}^p 3^{k-1} (\alpha_k \nu_k + \beta_k) \right), \quad \text{where } \Psi = \prod_{k=1}^p \psi_{k-1}. \end{aligned} \quad (41)$$

□

From (41) the stability of `solvegen` appears to depend on the parameters ψ_k , α_k and β_k . Unlike parameters α_k and β_k , parameter ψ_k can be easily estimated during the computation. It measures the magnitude of the matrices $F_r^{(k)}$ through the norms of the generators and can be used to monitor the behavior of the error. The parameter β_k has been introduced in the analysis of the errors arising in the computation of the generators of the matrices $F_r^{(k)}$. We expect it to be related to the conditioning of the matrices $A_r^{(k)}$ involved during the computation. If the matrix A has some special structural properties, the conditioning of all the matrices $A_r^{(k)}$ can be upper bounded by the conditioning of A . In particular, this happens if A is positive definite. Moreover, when A is diagonally dominant the norms $\|F_r^{(k)}\|$ are bounded and we expect that if A is well conditioned all the ψ_k are upper bounded.

Limiting our analysis to classes of matrices for which the parameters α_k , β_k and ψ_k are bounded by quantities $\hat{\alpha}$, $\hat{\beta}$ and $\hat{\psi}$ independent from N , we have

$$\pi_0 \leq \hat{\psi}^p \left(3^p \pi_e + \hat{\alpha} \sum_{k=1}^p 3^{k-1} n_k \log_2 n_k + \hat{\beta} \sum_{k=1}^p 3^{k-1} \right), \quad \text{where } n_k = 2^{p-k} n_e.$$

Then

$$\pi_0 \leq c_3 N^m, \quad \text{where } m = \log_2(3\hat{\psi}),$$

showing that π_0 depends at most polynomially on N .

5 Numerical experiments

The experiments have been conducted on an Intel Core Duo @ 3 GHz, 2GB RAM, using double precision arithmetic. Three sets of matrices have been considered:

(M1) The first set consists of 250 Toeplitz-like matrices of size 2^8 with displacement rank $\rho = 5$ and elementary size $n_e = 8$. These matrices are randomly generated with positive diagonal by means of the generators $C = D$ and of the signature matrix Σ , in such a way to give different conditioning and different values of the quantity $d = \min_i \left(1 - \sum_{j \neq i} |a_{i,j}|/|a_{ii}|\right)$ which is related to the point diagonal dominance of A . If $d > 0$, the matrix is diagonally dominant and positive definite, but positive definite matrices may occur also for negative values of d .

(M2) The second set consists of 250 symmetric positive definite Toeplitz matrices (displacement rank $\rho = 2$) of size 2^7 with elementary size $n_e = 4$. These matrices are generated as described in [13], through Schur parameters ρ_j , $j = 1, \dots, N$, distributed uniformly over symmetric intervals $\mathcal{I}_\ell = [-\ell, \ell]$, with $0.01 \leq \ell \leq 0.5$, and ρ_{10} and ρ_{15} changed to values with module close to one. The spectral properties of these matrices depend on ℓ and on the closeness to 1 of $|\rho_{10}|$ and $|\rho_{15}|$.

(M3) The 210 matrices of M3 are still generated through Schur parameters ρ_j , $j = 1, \dots, N$, distributed uniformly over the intervals \mathcal{I}_ℓ , with $0.001 \leq \ell \leq 0.35$, without further changes. When the endpoints of the intervals \mathcal{I}_ℓ are close to zero the generated matrices are diagonally dominant.

The left-hand side vector \mathbf{b} for the sets of matrices is computed from an exact solution, with elements randomly generated in $[0, 1]$. The considerations of the previous section suggest that the difficulty of solving a Toeplitz-like system by `solvegen` can be related to the conditioning of the matrices $A_r^{(k)}$ for any k . Since in general, the computation of the conditioning of all the matrices is unfeasible, we associate to each problem the parameter

$$\kappa = \max_r \mathbf{cond} (A_r^{(p)}),$$

where `cond` is the condition number. This parameter measures the greatest conditioning among the elementary blocks. In our experiments, the value of κ varies between $10^{0.9}$ and $10^{5.2}$ for M1 matrices and between $10^{0.5}$ and $10^{2.5}$ for M2 matrices. For M3 matrices the value of κ is bounded by $10^{1.5}$.

The scaled-independent relative residual

$$R = \frac{\|\mathbf{b} - A\tilde{\mathbf{x}}\|}{\|A\|\|\tilde{\mathbf{x}}\| + \|\mathbf{b}\|}$$

can be taken as an experimental estimate of the perturbation $\|\Delta_A\|$ occurring in Theorem 3. Of course, R is expected to have a better behavior than the bound

appearing in the theorem. The aim of the experimentation is to point out that the quantity Ψ , which appears in the theoretical bound (41) of the perturbation $\|\Delta_A\|$ plays a relevant role also in practice.

Figure 1 shows the residuals R (black points) plotted versus κ in Log-Log scale for problems M1 on the left and M2 on the right. The figure shows also the corresponding values of Ψ (gray points), suitably scaled by a factor, in order to put the graphs in the same figure. It appears that parameter Ψ well describes the behavior of the residual. Moreover the points lie in a strip around a straight line, pointing out a polynomial growth of R and Ψ with respect to κ .

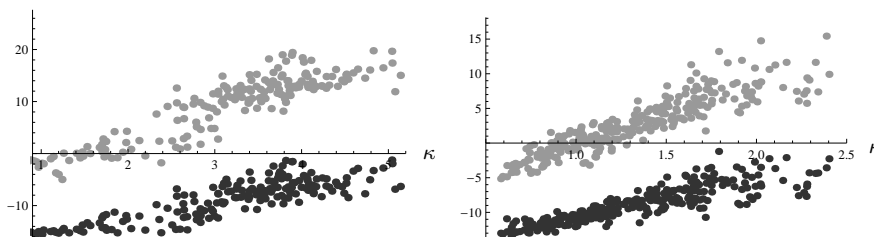


Figure 1: Log-Log plot of the relative residuals R (black points) and of $\Psi/10^{10}$ (gray points) as functions of κ , for problems M1 (on the left) and problems M2 (on the right).

Figure 2 shows the residual in Log scale plotted versus d for matrices M3. More precisely, the figure on the left refers to matrices generated through ρ_j belonging to \mathcal{I}_ℓ , with $0.01 \leq \ell \leq 0.35$. These matrices, even if not diagonally dominant, are well conditioned with small absolute values of parameter d . The figure on the right refers to diagonally dominant matrices generated through ρ_j belonging to \mathcal{I}_ℓ , with $0.001 \leq \ell \leq 0.015$. It is evident that for these matrices `solvegen` has a good performance, the stronger the diagonally dominance, the better the performance. From the figure on the left it appears that parameter Ψ may result in a pessimistic estimate of the residual especially for problems with good spectral properties.

6 Conclusions

In this paper we have performed a theoretical error analysis of a superfast algorithm recently introduced by the authors in [4]. This analysis allows the characterization of some parameters which rule the stability behavior of the algorithm. Among them, the parameters ψ_k , measuring the magnitude of the matrices $F_k^{(r)}$ through their generators, can be easily estimated during the computation. The numerical experimentation points out that the quantity $\Psi = \prod_{k=1}^p \psi_{k-1}$, which

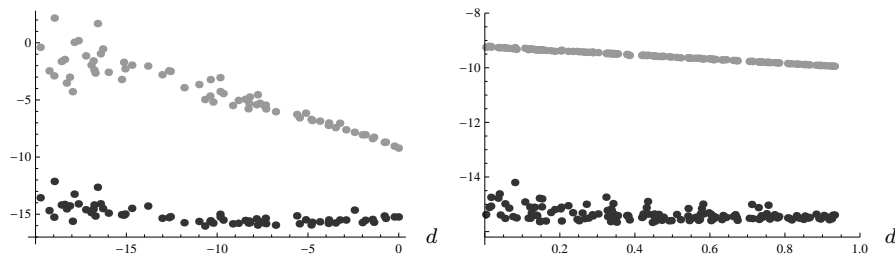


Figure 2: Log plot of the relative residuals R (black points) and of $\Psi/10^{10}$ (gray points) as functions of d , for problems M2.

appears in the theoretical bound (41), well describes also in practice the residual behavior. Anyway, the use of Ψ for estimating the residual may be too pessimistic, especially when `solvegen` has a good performance.

References

- [1] M. Arioli, H. Munthe-Kaas, L. Valdettaro, Componentwise error analysis for FFT's with applications to fast Helmholtz solvers, *Numer. Algorithms*, 12, (1996), pp. 65-88.
- [2] R. R. Bitmead and B. D. O. Anderson, "Asymptotically Fast Solution of Toeplitz and Related Systems of Linear Equations", *Linear Algebra Appl.*, 34, pp. 103-116, 1980.
- [3] P. Favati, G. Lotti and O. Menchi, "Stability of the Levinson algorithm for Toeplitz-like systems", *SIAM Journal on Matrix Analysis and Applications*, 31, pp. 2531-2552, 2010.
- [4] P. Favati, G. Lotti and O. Menchi, "Superfast solution of Toeplitz-like systems", Tech. Report IIT TR-24/2011.
- [5] P. Favati, G. Lotti and O. Menchi, "Stability analysis of the product via FFT of a Toeplitz-like matrix", Tech. Report IIT TR-10/2014.
- [6] M. Gu, "Stable and Efficient Algorithms for Structured Systems of Linear Equations", *SIAM Journal on Matrix Analysis and Applications*, 19, pp. 279-306, 1998.
- [7] G. Heinig and K. Rost, *Algebraic methods for Toeplitz-like matrices and operators*, Akademie-Verlag, Berlin, 1984.
- [8] N. J. Higham, *Accuracy and Stability of Numerical Algorithms*, SIAM, Philadelphia, PA, 1996.
- [9] T. Kailath, S.-Y. Kung and M. Morf, "Displacement ranks of matrices and linear equations", *J. Math. Anal. Appl.*, 68, pp. 395-407, 1979.

- [10] T. Kailath and A. H. Sayed, "Displacement structure: theory and applications", *SIAM Rev.*, 37, pp. 297-386, 1995.
- [11] T. Kailath, A. Viera, and M. Morf, Inverses of Toeplitz operators, innovations and orthogonal polynomials, *SIAM Rev.*, 20, pp. 106-119, 1978.
- [12] M. Morf, Doubling algorithms for Toeplitz and related equations, in *Proc. IEEE Internat. conf. on Acoustics, Speech, & Signal Processing*, pp. 954-959, Denver, Colorado, 1980.
- [13] M. Stewart, "A Superfast Toeplitz Solver with Improved Numerical Stability", *SIAM J. Matrix Anal. Appl.*, 25, pp. 669-693, 2003.