

# ON CIRCULANT AND SKEW-CIRCULANT SPLITTING ALGORITHMS FOR (CONTINUOUS) SYLVESTER EQUATIONS \*

ZHONGYUN LIU<sup>†</sup>, FANG ZHANG<sup>\*</sup>, CARLA FERREIRA<sup>‡</sup>, AND YULIN ZHANG<sup>‡</sup>

**Abstract.** We present a circulant and skew-circulant splitting (CSCS) iterative method for solving large sparse continuous Sylvester equations  $AX + XB = C$ , where the coefficient matrices  $A$  and  $B$  are Toeplitz matrices. A theoretical study shows that if the circulant and skew-circulant splitting factors of  $A$  and  $B$  are positive semi-definite and at least one is positive definite (not necessarily Hermitian), then the CSCS method converges to the unique solution of the Sylvester equation. In addition, we obtain an upper bound for the convergence factor of the CSCS iteration. This convergence factor depends only on the eigenvalues of the circulant and skew-circulant splitting matrices. A computational comparison with alternative methods reveals the efficiency and reliability of the proposed method.

**Key words.** Continuous Sylvester equations, CSCS iteration, Toeplitz matrices, convergence

**AMS subject classifications.** 15A24, 65F10, 65H10

**1. Introduction.** A continuous Sylvester equation is possibly one of the most popular linear matrix equations used in mathematics. It is a matrix equation of the form

$$AX + XB = C, \tag{1.1}$$

where matrices  $A \in \mathbb{C}^{n \times n}$ ,  $B \in \mathbb{C}^{m \times m}$ ,  $C \in \mathbb{C}^{n \times m}$  are given and the problem is to find a matrix  $X \in \mathbb{C}^{n \times m}$  that obeys this equation. It is well-known that equation (1.1) has a unique solution for  $X$  if and only if  $A$  and  $-B$  do not have common eigenvalues (see, e.g., [23, 28]).

The Sylvester equation is classically employed in the design of Luenberger observers, which are widely used in signal processing, control and system theory (see, e.g., [8, 10, 16, 20, 27]); often appears in linear and generalized eigenvalue problems for the Riccati equation in the computation of invariant subspaces (see, e.g., [6, 17, 38]); can be used to devise implicit Runge-Kutta integral formulae and block multi-step formulae for the numerical solutions of ordinary differential equations (see, e.g., [19]); and some linear systems arising, for example, from finite difference discretizations of separable elliptic boundary value problems on rectangular domains, can be written as a Sylvester equation (see, e.g., [14, 18]).

There are essentially two different approaches to deal with the Sylvester equation (1.1). The first approach consists in vectorizing the unknown matrix  $X$  and translating the matrix equation into a linear system  $\mathcal{A}\mathbf{x} = \mathbf{c}$ , where vectors  $\mathbf{x}$  and  $\mathbf{c}$  are the column-stacking vectors of the matrices  $X$  and  $C$ , respectively, and  $\mathcal{A}$  is the *Kronecker sum* of the matrices  $A$  and  $B^T$ , that is,  $\mathcal{A} = I_m \otimes A + B^T \otimes I_n$ , with symbol  $\otimes$  denoting the standard Kronecker product. Either direct or iterative methods can be applied to solve this linear system. The second approach is to treat the Sylvester equation (1.1) in its original form using an iterative method directly applied to matrices  $A$ ,  $B$  and  $C$ .

When matrices  $A$  and  $B$  are large, the order of the coefficient matrix  $\mathcal{A} \in \mathbb{C}^{mn \times mn}$  in the linear system

---

\* The research of the last two authors was partially financed by Portuguese Funds through FCT (Fundação para a Ciência e a Tecnologia) within the Projects UIDB/00013/2020 and UIDP/00013/2020. This work was also supported by NSFC (National Natural Science Foundation of China) under the Grant No. 11371075, the Hunan Key Laboratory of Mathematical Modeling and Analysis in Engineering.

<sup>†</sup>School of Mathematics and Statistics, Changsha University of Science and Technology, Changsha 410076, P. R. China (liuzhongyun@263.net).

<sup>‡</sup>Centro de Matemática, Universidade do Minho, 4710-057 Braga, Portugal (caferrei@math.uminho.pt, zhang@math.uminho.pt).

$\mathcal{A}\mathbf{x} = \mathbf{c}$  will be considerably larger and, in general, difficulties related to data storage and computational time arise. This explains that the first approach is mainly used in problems of small or medium dimension. The Bartels-Stewart method proposed in [7] is based on the reduction of the matrices  $A$  and  $B$  to real Schur form (quasi-triangular form) using the QR algorithm for eigenvalues, followed by the use of direct methods to solve several linear systems. The Hessenberg-Schur method, presented in [21], reduces matrix  $A$  to Hessenberg form and only matrix  $B$  is decomposed into the quasi-triangular Shur form and it is faster than the Bartels-Stewart method. However, in both methods, the authors were unable to establish a backward stability result. These methods are classified as direct methods and are used by MATLAB.

When matrices  $A$  and  $B$  are large and sparse, following the second approach, iterative methods such as the Smith's method [37], the *alternating direction implicit* method (ADI) [9, 12, 24, 31, 40], the *block successive over-relaxation* method (BSOR) [35] and the *matrix splitting* methods [1, 22] are efficient and accurate methods to obtain a numerical solution of the equation (1.1). The development of the mentioned iterative methods based on the concept of matrix splitting has attracted several scholars and a large number of efficient and robust algorithms were proposed. See [2, 29, 30, 41, 44] and references therein.

In this paper, we consider the case when  $A$  and  $B$  are both Toeplitz matrices. Matrices with this structure appear, for example, in connection to the discretization of the convection-diffusion reaction equation [14, 18]. We present an iterative method for solving the Sylvester equation (1.1) using the circulant and skew-circulant splittings of the matrices  $A$  and  $B$ . This circulant and skew-circulant splitting (CSCS) iteration method is a matrix variant of the CSCS iteration method firstly proposed in [33] for solving a Toeplitz linear system. These type of methods are conceptually analogous to the ADI iteration methods. Via this CSCS iteration method, the problem of solving a general continuous Sylvester equation is translated into two coupled continuous Sylvester equations involving shifted circulant and skew-circulant matrices.

When the circulant and skew-circulant splitting matrices of  $A$  and  $B$  are positive definite (not necessarily Hermitian), we prove that the CSCS iteration converges unconditionally to the exact solution of the Sylvester equation (1.1). Moreover, the values of the shift parameters that minimize an upper bound for the contraction factor are obtained in terms of the bounds for the largest and the smallest eigenvalues of the circulant and skew-circulant splitting matrices of  $A$  and  $B$ .

The organization of this paper is as follows. After giving some basic definitions and preliminary results in section 2, we describe the CSCS iterative method for solving equation (1.1) in section 3. We then analyze some sufficient conditions that ensure the convergence of this method in section 4. Numerical experiments are shown in section 5. These examples illustrate the efficiency and robustness of our method.

**2. Basic definitions and preliminary results.** Given a matrix  $K \in \mathbb{C}^{n \times m}$ ,  $K^*$  denotes the conjugate transpose of  $K$ , the  $(i, j)$  element of  $K$  is denoted by  $K_{i,j}$  and  $\rho(K)$  stands for the spectral radius of  $K$ . The set of all the eigenvalues of  $K$  is represented by  $\lambda(K)$ . If  $x \in \mathbb{C}$ ,  $\text{Re}(x)$  denotes the real part of  $x$  and  $\text{Im}(x)$  the imaginary part.

Here we use the general concept of positive definiteness which says that a matrix  $K \in \mathbb{C}^{n \times n}$  is positive definite if its Hermitian part  $\frac{1}{2}(K + K^*)$  is positive definite in the narrower sense. In general, this condition is equivalent to  $\text{Re}(z^* K z) > 0$ , for all nonzero vectors  $z \in \mathbb{C}$ , which implies that  $\text{Re}(\lambda) > 0$ , for any eigenvalue  $\lambda$  of  $K$ .

A square matrix  $T$  is said to be Toeplitz (or diagonal-constant) when  $T_{j,k} = t_{j-k}$ ,  $j, k = 1, \dots, n$ , for constants  $t_{1-n}, \dots, t_{n-1}$ . An important property of a Toeplitz matrix  $T$  is that it always admits the additive decomposition

$$T = C_T + S_T, \tag{2.1}$$

where  $C_T$  is a circulant matrix and  $S_T$  is a skew-circulant matrix. See [33, 34]. We say that a matrix  $C$  is

a circulant matrix if  $C_{j,k} = c_{j-k}$ , for constants  $c_{1-n}, \dots, c_{n-1}$  such that  $c_{-l} = c_{n-l}$ ,  $l = 1, \dots, n-1$ . That is, a circulant matrix is a Toeplitz matrix that is fully defined by its first column (or row) given that the remaining columns are cyclic permutations of the first column (or row). A skew-circulant matrix is also a particular type of Toeplitz matrix. We say that a matrix  $S$  is a skew-circulant matrix if  $S_{j,k} = s_{j-k}$ , for constants  $s_{1-n}, \dots, s_{n-1}$  such that  $s_{-l} = -s_{n-l}$ ,  $l = 1, \dots, n-1$ .

Matrices  $C_T$  and  $S_T$  in (2.1), the circulant and skew-circulant splitting (CSCS) of  $T$ , are defined as follows:

$$(C_T)_{j,k} = \frac{1}{2} \begin{cases} t_0, & \text{if } j = k, \\ t_{j-k} + t_{j-k+n}, & \text{if } j < k, \\ t_{j-k} + t_{j-k-n}, & \text{if } j > k, \end{cases} \quad \text{and} \quad (S_T)_{j,k} = \frac{1}{2} \begin{cases} t_0, & \text{if } j = k, \\ t_{j-k} - t_{j-k+n}, & \text{if } j < k, \\ t_{j-k} - t_{j-k-n}, & \text{if } j > k. \end{cases} \quad (2.2)$$

It is well-known that a circulant matrix  $C$  is diagonalizable by the unitary Fourier matrix  $F$  of order  $n$  which entries are given by

$$F_{j,k} = \frac{1}{\sqrt{n}} \omega^{(j-1)(k-1)}, \quad j, k = 1, \dots, n,$$

where  $\omega$  is the primitive  $n$ -th root of unit  $\omega = e^{\frac{2\pi}{n}\mathbf{i}}$ ,  $\mathbf{i} = \sqrt{-1}$ . Similarly, a skew-circulant matrix  $S$  is diagonalizable by the unitary matrix  $\hat{F} = FD$  where  $D = \text{diag}(1, e^{\frac{\pi}{n}\mathbf{i}}, \dots, e^{\frac{(n-1)\pi}{n}\mathbf{i}})$ . Thus,

$$C = F^* \Lambda F \quad \text{and} \quad S = \hat{F}^* \Sigma \hat{F}, \quad (2.3)$$

where  $\Lambda$  and  $\Sigma$  are diagonal matrices holding the eigenvalues of  $C$  and  $S$ , respectively. Moreover, we note that  $\Lambda$  and  $\Sigma$  can be obtained in  $\mathcal{O}(n \log n)$  operations by taking the fast Fourier transform (FFT) of the first column (or row) of  $C$  and first row of  $S$ , respectively. In fact, the diagonal entries  $\lambda_j$  of  $\Lambda$  and the diagonal entries  $\sigma_j$  of  $\Sigma$  are given, respectively, by

$$\lambda_j = \sum_{k=1}^n c_{k-1} \omega^{(j-1)(k-1)} \quad \text{and} \quad \sigma_j = \sum_{k=1}^n s_{k-1} \omega^{(j-1)(k-1)} e^{\frac{\pi(k-1)\pi}{n}\mathbf{i}}, \quad j = 1, \dots, n. \quad (2.4)$$

See, for instance, [13, 15]. For the FFT algorithm, we refer to [32, 39].

Once  $\Lambda$  and  $\Sigma$  are obtained, the products  $C\mathbf{x}$  and  $C^{-1}\mathbf{x}$ , as well as  $S\mathbf{x}$  and  $S^{-1}\mathbf{x}$ , for any vector  $\mathbf{x}$ , can be computed by FFTs in  $\mathcal{O}(n \log n)$  operations. Therefore, the use of circulant and skew-circulant matrices to solve matrix equations with Toeplitz matrices allows to improve the efficiency by employing FFTs throughout the computations. For a matrix  $M \in \mathbb{C}^{n \times m}$ , the FFT operation is applied to each column in  $\mathcal{O}(mn \log n)$ .

**3. CSCS Iteration.** Let matrices  $A \in \mathbb{C}^{n \times n}$  and  $B \in \mathbb{C}^{m \times m}$  have a Toeplitz structure and let

$$A = C_A + S_A, \quad B = C_B + S_B \quad (3.1)$$

be the circulant and skew-circulant splittings of  $A$  and  $B$  (CSCS), respectively. There are no constraints in using (2.2), so these splittings always exist.

If  $\alpha$  and  $\beta$  are positive constants, the following splittings are also CSCS splittings of  $A$  and  $B$ ,

$$A = (C_A + \alpha I_n) + (S_A - \alpha I_n), \quad B = (C_B + \beta I_m) + (S_B - \beta I_m), \quad (3.2)$$

$$A = (C_A - \alpha I_n) + (S_A + \alpha I_n), \quad B = (C_B - \beta I_m) + (S_B + \beta I_m). \quad (3.3)$$

It follows that if  $X^* \in \mathbb{C}^{n \times m}$  is the exact solution of the Sylvester equation (1.1), then

$$\begin{cases} (C_A + \alpha I)X^* + X^*(C_B + \beta I) = (\alpha I - S_A)X^* + X^*(\beta I - S_B) + C \\ (S_A + \alpha I)X^* + X^*(S_B + \beta I) = (\alpha I - C_A)X^* + X^*(\beta I - C_B) + C \end{cases}$$

where  $I$  is either the identity matrix of order  $n$  or  $m$ , conformable to  $C_A$  or  $C_B$ , respectively. We are now able to define the fixed-point matrix equations

$$\begin{cases} (C_A + \alpha I)X + X(C_B + \beta I) = (\alpha I - S_A)Y + Y(\beta I - S_B) + C \\ (S_A + \alpha I)Y + Y(S_B + \beta I) = (\alpha I - C_A)X + X(\beta I - C_B) + C \end{cases} \quad (3.4)$$

such that  $X^*$  is the fixed point of both equations. The reverse also occurs, that is, if  $X^*$  is a fixed point of either of the two equations in (3.4), then it is the exact solution of (1.1). See [4, Theorem 3.1].

The CSCS iteration method is defined as follows.

**CSCS iteration method.** *Given an initial approximation  $X^{(0)}$  and positive constants  $\alpha, \beta$  (shift parameters), repeat the iterative scheme*

$$\begin{cases} (\alpha I + C_A)X^{(k+\frac{1}{2})} + X^{(k+\frac{1}{2})}(\beta I + C_B) = (\alpha I - S_A)X^{(k)} + X^{(k)}(\beta I - S_B) + C & \left( \text{solve for } X^{(k+\frac{1}{2})} \right) \\ (\alpha I + S_A)X^{(k+1)} + X^{(k+1)}(\beta I + S_B) = (\alpha I - C_A)X^{(k+\frac{1}{2})} + X^{(k+\frac{1}{2})}(\beta I - C_B) + C & \left( \text{solve for } X^{(k+1)} \right) \end{cases} \quad (3.5)$$

for  $k = 0, 1, 2, \dots$ , until  $\{X^{(k)}\}$  converges.

An alternative version of this two-step iteration is obtained if we compute the corrections  $Z^{(k+\frac{1}{2})} = X^{(k+\frac{1}{2})} - X^{(k)}$  and  $Z^{(k+1)} = X^{(k+1)} - X^{(k+\frac{1}{2})}$  in each iteration which brings the residuals  $R^{(k)}$  and  $R^{(k+\frac{1}{2})}$  into the computation. The iterative scheme is changed to

$$\begin{cases} R^{(k)} = C - AX^{(k)} - X^{(k)}B \\ (\alpha I + C_A)Z^{(k+\frac{1}{2})} + Z^{(k+\frac{1}{2})}(\beta I + C_B) = R^{(k)} & \left( \text{solve for } Z^{(k+\frac{1}{2})} \right) \\ X^{(k+\frac{1}{2})} = X^{(k)} + Z^{(k+\frac{1}{2})} \\ R^{(k+\frac{1}{2})} = C - AX^{(k+\frac{1}{2})} - X^{(k+\frac{1}{2})}B \\ (\alpha I + S_A)Z^{(k+1)} + Z^{(k+1)}(\beta I + S_B) = R^{(k+\frac{1}{2})} & \left( \text{solve for } Z^{(k+1)} \right) \\ X^{(k+1)} = X^{(k+\frac{1}{2})} + Z^{(k+1)} \end{cases} \quad (3.6)$$

Each CSCS iteration requires the solution of two Sylvester equations. In the first step  $Z^{(k+\frac{1}{2})}$  is the solution of the equation

$$(\alpha I + C_A)Z + Z(\beta I + C_B) = R^{(k)} \quad (3.7)$$

and in the second step  $Z^{(k+1)}$  is the solution of the equation

$$(\alpha I + S_A)Z + Z(\beta I + S_B) = R^{(k+\frac{1}{2})}. \quad (3.8)$$

The method alternates between equation (3.7), with circulant matrices  $C_A$  and  $C_B$ , and equation (3.8), with skew-circulant matrices  $S_A$  and  $S_B$ , and we can reverse the roles of these matrix equations.

Observe that once we compute the eigenvalues of  $C_A$  and  $C_B$ , it is always possible to choose positive constants  $\alpha$  and  $\beta$  such that  $\alpha I + C_A$  and  $-(\beta I + C_B)$  do not have eigenvalues in common and thus the matrix equation (3.7) has a unique solution. A similar observation applies to the matrix equation (3.8) concerning the matrices  $\alpha I + S_A$  and  $-(\beta I + S_B)$ .

Obtained the eigenvalue decompositions, see (2.3),

$$C_A = F_A^* \Lambda_A F_A, \quad C_B = F_B^* \Lambda_B F_B, \quad S_A = \hat{F}_A^* \Sigma_A \hat{F}_A, \quad S_B = \hat{F}_B^* \Sigma_B \hat{F}_B, \quad (3.9)$$

equation (3.7) is equivalent to

$$(\alpha I + \Lambda_A)\mathcal{Z} + \mathcal{Z}(\beta I + \Lambda_B) = \mathcal{R}^{(k)} \quad (3.10)$$

where  $\mathcal{Z} = F_A Z F_B^*$  and  $\mathcal{R}^{(k)} = F_A R^{(k)} F_B^*$ ; and equation (3.8) is equivalent to

$$(\alpha I + \Sigma_A) \mathcal{Z} + \mathcal{Z} (\beta I + \Sigma_B) = \mathcal{R}^{(k+\frac{1}{2})} \quad (3.11)$$

where  $\mathcal{Z} = \hat{F}_A Z \hat{F}_B^*$  and  $\mathcal{R}^{(k+\frac{1}{2})} = \hat{F}_A R^{(k+\frac{1}{2})} \hat{F}_B^*$ .

Solving the Sylvester equations (3.10) and (3.11) is immediate since these matrix equations can be translated into linear systems with diagonal coefficient matrices,  $I_m \otimes (\alpha I_n + \Lambda_A) + (\beta I_m + \Lambda_B)^T \otimes I_n$  and  $I_m \otimes (\alpha I_n + \Sigma_A) + (\beta I_m + \Sigma_B)^T \otimes I_n$ , respectively.

The solutions of the initial Sylvester equations (3.7) and (3.8) are given by  $Z = F_A^* \mathcal{Z} F_B$ , for  $\mathcal{Z}$  satisfying (3.10), and  $Z = \hat{F}_A^* \mathcal{Z} \hat{F}_B$ , for  $\mathcal{Z}$  satisfying (3.11), respectively. These products can be computed efficiently using FFTs.

It is also possible, once again using FFTs, to reduce the computational effort associated to the right-hand side of the two Sylvester equations involved in each CSCS iteration, equations (3.10) and (3.11). For an approximation  $X^{(j)}$ , the residual  $R^{(j)}$  can be computed using the decomposition

$$\begin{aligned} R^{(j)} &= C - (C_A + S_A) X^{(j)} - X^{(j)} (C_B + S_B) \\ &= C - F_A^* \Lambda_A F_A X^{(j)} - \hat{F}_A^* \Sigma_A \hat{F}_A X^{(j)} - X^{(j)} F_B^* \Lambda_B F_B - X^{(j)} \hat{F}_B^* \Sigma_B \hat{F}_B. \end{aligned} \quad (3.12)$$

Thus, the right-hand sides of equations (3.10) and (3.11) are given, respectively, by

$$\mathcal{R}^{(k)} = F_A \left[ C - F_A^* \Lambda_A F_A X^{(k)} - \hat{F}_A^* \Sigma_A \hat{F}_A X^{(k)} - X^{(k)} F_B^* \Lambda_B F_B - X^{(k)} \hat{F}_B^* \Sigma_B \hat{F}_B \right] F_B^* \quad (3.13)$$

and

$$\mathcal{R}^{(k+\frac{1}{2})} = \hat{F}_A \left[ C - F_A^* \Lambda_A F_A X^{(k+\frac{1}{2})} - \hat{F}_A^* \Sigma_A \hat{F}_A X^{(k+\frac{1}{2})} - X^{(k+\frac{1}{2})} F_B^* \Lambda_B F_B - X^{(k+\frac{1}{2})} \hat{F}_B^* \Sigma_B \hat{F}_B \right] \hat{F}_B^*. \quad (3.14)$$

Given an initial approximation  $X^{(0)}$ , positive constants  $\alpha$ ,  $\beta$ , and the spectral factorizations (3.9), the two steps of the CSCS iteration (3.6) can then be expressed as

$$\left\{ \begin{array}{l} \text{Use (3.13) to compute } \mathcal{R}^{(k)} \\ (\alpha I + \Lambda_A) \mathcal{Z} + \mathcal{Z} (\beta I + \Lambda_B) = \mathcal{R}^{(k)} \quad \text{(solve for } \mathcal{Z}) \\ X^{(k+\frac{1}{2})} = X^{(k)} + F_A^* \mathcal{Z} F_B \\ \text{Use (3.14) to compute } \mathcal{R}^{(k+\frac{1}{2})} \\ (\alpha I + \Sigma_A) \mathcal{Z} + \mathcal{Z} (\beta I + \Sigma_B) = \mathcal{R}^{(k+\frac{1}{2})} \quad \text{(solve for } \mathcal{Z}) \\ X^{(k+1)} = X^{(k+\frac{1}{2})} + \hat{F}_A^* \mathcal{Z} \hat{F}_B \end{array} \right. \quad (3.15)$$

for  $k = 0, 1, 2, \dots$ , until  $\{X^{(k)}\}$  converges.

Notice that all the matrix multiplications can be performed using FFTs and thus the operation count is  $\mathcal{O}(mn \log n)$  (or  $\mathcal{O}(nm \log m)$ , depending on which is bigger). There is no need to perform explicit matrix multiplications. See Algorithm 1 in Appendix A for a MATLAB implementation of CSCS.

**4. Convergence results.** Using a matrix-vector formulation of the Sylvester equation (1.1), such that vectors  $\mathbf{x}$  and  $\mathbf{c}$  are the column-stacking vectors of the matrices  $X$  and  $C$ , respectively, the two-step iterative CSCS scheme (3.5) can be rewritten as

$$\begin{cases} [I_m \otimes (\alpha I + C_A) + (\beta I + C_B)^T \otimes I_n] \mathbf{x}^{(k+\frac{1}{2})} = [I_m \otimes (\alpha I - S_A) + (\beta I - S_B)^T \otimes I_n] \mathbf{x}^{(k)} + \mathbf{c} \\ [I_m \otimes (\alpha I + S_A) + (\beta I + S_B)^T \otimes I_n] \mathbf{x}^{(k+1)} = [I_m \otimes (\alpha I - C_A) + (\beta I - C_B)^T \otimes I_n] \mathbf{x}^{(k+\frac{1}{2})} + \mathbf{c} \end{cases} \quad (4.1)$$

for  $k = 0, 1, 2, \dots$ , until  $\{\mathbf{x}^{(k)}\}$  converges, given an initial approximation  $\mathbf{x}^{(0)}$  and positive constants  $\alpha, \beta$ .

In this section we will establish theoretical results concerning the convergence conditions of the CSCS iteration and our analysis is based on this matrix-vector formulation of the method.

So, we are considering the linear system  $\mathcal{A}\mathbf{x} = \mathbf{c}$ , equivalent to the Sylvester equation (1.1), where  $\mathcal{A} = I_m \otimes A + B^T \otimes I_n$ , and the two splittings of the matrix  $\mathcal{A}$ ,

$$\begin{aligned}\mathcal{A} &= [I_m \otimes (\alpha I + C_A) + (\beta I + C_B)^T \otimes I_n] - [I_m \otimes (\alpha I - S_A) + (\beta I - S_B)^T \otimes I_n], \\ \mathcal{A} &= [I_m \otimes (\alpha I + S_A) + (\beta I + S_B)^T \otimes I_n] - [I_m \otimes (\alpha I - C_A) + (\beta I - C_B)^T \otimes I_n],\end{aligned}$$

corresponding to the CSCS splittings of  $A$  and  $B$  given in (3.2). Using the bilinearity property of the kronecker product, these decompositions of  $\mathcal{A}$  can be rewritten as

$$\begin{aligned}\mathcal{A} &= [(\alpha + \beta)I_{mn} + (I_m \otimes C_A + C_B^T \otimes I_n)] - [(\alpha + \beta)I_{mn} - (I_m \otimes S_A + S_B^T \otimes I_n)], \\ \mathcal{A} &= [(\alpha + \beta)I_{mn} + (I_m \otimes S_A + S_B^T \otimes I_n)] - [(\alpha + \beta)I_{mn} - (I_m \otimes C_A + C_B^T \otimes I_n)].\end{aligned}$$

Thereby, defining  $\tilde{C} = (I_m \otimes C_A + C_B^T \otimes I_n)$  and  $\tilde{S} = (I_m \otimes S_A + S_B^T \otimes I_n)$ , we have

$$\begin{aligned}\mathcal{A} &= [(\alpha + \beta)I + \tilde{C}] - [(\alpha + \beta)I - \tilde{S}], \\ \mathcal{A} &= [(\alpha + \beta)I + \tilde{S}] - [(\alpha + \beta)I - \tilde{C}],\end{aligned}\tag{4.2}$$

where  $I$  represents the identity matrix of order  $mn$ , and the iteration (4.1) can be expressed as

$$\begin{cases} [(\alpha + \beta)I + \tilde{C}]\mathbf{x}^{(k+\frac{1}{2})} = [(\alpha + \beta)I - \tilde{S}]\mathbf{x}^{(k)} + \mathbf{c} \\ [(\alpha + \beta)I + \tilde{S}]\mathbf{x}^{(k+1)} = [(\alpha + \beta)I - \tilde{C}]\mathbf{x}^{(k+\frac{1}{2})} + \mathbf{c}. \end{cases}\tag{4.3}$$

This iterative scheme is a particular case of a more general two-step splitting scheme defined by

$$\begin{cases} M_1\mathbf{x}^{(k+\frac{1}{2})} = N_1\mathbf{x}^{(k)} + \mathbf{c} \\ M_2\mathbf{x}^{(k+1)} = N_2\mathbf{x}^{(k+\frac{1}{2})} + \mathbf{c}, \quad k = 0, 1, 2, \dots, \end{cases}\tag{4.4}$$

assuming that  $\mathcal{A}$  admits the decompositions  $\mathcal{A} = M_i - N_i$ ,  $i = 1, 2$ , with  $M_1$  and  $M_2$  invertible matrices. The sequence  $\{\mathbf{x}^{(k)}\}$  generated by (4.4) satisfies

$$\mathbf{x}^{(k+1)} = \mathcal{M}\mathbf{x}^{(k)} + \mathcal{C}\mathbf{c}, \quad k = 0, 1, 2, \dots,\tag{4.5}$$

where

$$\mathcal{M} = M_2^{-1}N_2M_1^{-1}N_1 \quad \text{and} \quad \mathcal{C} = M_2^{-1}(I + N_2M_1^{-1}).\tag{4.6}$$

Matrix  $\mathcal{M}$  is called the *iteration matrix* and it is well-known that  $\{\mathbf{x}^{(k)}\}$  converges to the exact solution of the linear system  $\mathcal{A}\mathbf{x} = \mathbf{c}$  if and only if  $\rho(\mathcal{M}) < 1$ , for any initial approximation  $\mathbf{x}^{(0)} \in \mathbb{C}^{mn}$  [36, 42].

We will prove that when  $\tilde{C} = I_m \otimes C_A + C_B^T \otimes I_n$  and  $\tilde{S} = I_m \otimes S_A + S_B^T \otimes I_n$  are positive definite matrices (the real part of the eigenvalues is positive), then the proposed CSCS iterative method converges.

**THEOREM 4.1.** *Let  $A \in \mathbb{C}^{n \times n}$  and  $B \in \mathbb{C}^{m \times m}$  be Toeplitz matrices such that  $A = C_A + S_A$  and  $B = C_B + S_B$  are the circulant and skew-circulant splittings of  $A$  and  $B$ , respectively. Consider the linear system  $\mathcal{A}\mathbf{x} = \mathbf{c}$ , where  $\mathcal{A} = I_m \otimes A + B^T \otimes I_n$ , equivalent to the Sylvester equation (1.1), and the splitting  $\mathcal{A} = \tilde{C} + \tilde{S}$ , where*

$$\tilde{C} = I_m \otimes C_A + C_B^T \otimes I_n \quad \text{and} \quad \tilde{S} = I_m \otimes S_A + S_B^T \otimes I_n.\tag{4.7}$$

Let  $\alpha, \beta$  be two positive constants and  $\gamma = \alpha + \beta$ . Then the iteration matrix of the CSCS scheme (4.3) is

$$\mathcal{M}_\gamma = (\gamma I + \tilde{S})^{-1}(\gamma I - \tilde{C})(\gamma I + \tilde{C})^{-1}(\gamma I - \tilde{S}) \quad (4.8)$$

and its spectral radius  $\rho(\mathcal{M}_\gamma)$  is bounded by

$$\sigma_\gamma \equiv \max_{\lambda_j \in \lambda(\tilde{C})} \left| \frac{\gamma - \lambda_j}{\gamma + \lambda_j} \right| \cdot \max_{\mu_j \in \lambda(\tilde{S})} \left| \frac{\gamma - \mu_j}{\gamma + \mu_j} \right|.$$

If  $\tilde{C}$  is positive definite and  $\tilde{S}$  is positive semi-definite (or vice-versa), then

$$\rho(\mathcal{M}_\gamma) \leq \sigma_\gamma < 1, \quad \text{for all } \gamma > 0,$$

and, thus, the CSCS iteration (4.3) converges to the exact solution  $\mathbf{x}^*$  of the linear system  $\mathcal{A}\mathbf{x} = \mathbf{c}$ . The equivalent CSCS iteration (3.5) converges to the exact solution  $X^* \in \mathbb{C}^{m \times n}$  of the Sylvester equation (1.1).

*Proof.* Given  $\gamma = \alpha + \beta$ , the CSCS iteration (4.3) can be rewritten as

$$\begin{cases} (\gamma I + \tilde{C})\mathbf{x}^{(k+\frac{1}{2})} = (\gamma I - \tilde{S})\mathbf{x}^{(k)} + \mathbf{c} \\ (\gamma I + \tilde{S})\mathbf{x}^{(k+1)} = (\gamma I - \tilde{C})\mathbf{x}^{(k+\frac{1}{2})} + \mathbf{c} \end{cases} \quad (4.9)$$

which is the two-step splitting iterative scheme (4.4) to solve the linear system  $\mathcal{A}\mathbf{x} = \mathbf{c}$  with  $M_1 = \gamma I + \tilde{C}$ ,  $N_1 = \gamma I - \tilde{S}$ ,  $M_2 = \gamma I + \tilde{S}$  and  $N_2 = \gamma I - \tilde{C}$ . According to (4.5) and (4.6), the two steps of iteration (4.9) can be put together into the stationary fixed-point iteration

$$\mathbf{x}^{(k+1)} = \mathcal{M}_\gamma \mathbf{x}^{(k)} + \mathcal{E}_\gamma \mathbf{c}, \quad k = 0, 1, 2, \dots, \quad (4.10)$$

where

$$\mathcal{M}_\gamma = (\gamma I + \tilde{S})^{-1}(\gamma I - \tilde{C})(\gamma I + \tilde{C})^{-1}(\gamma I - \tilde{S}) \quad \text{and} \quad \mathcal{E}_\gamma = 2\gamma(\gamma I + \tilde{S})^{-1}(\gamma I + \tilde{C})^{-1}.$$

The spectral radius  $\rho(\mathcal{M}_\gamma)$  governs the convergence of (4.10) and, since the spectrum of a matrix is invariant under a similarity transformation, we find that

$$\begin{aligned} \rho(\mathcal{M}_\gamma) &= \rho\left((\gamma I - \tilde{C})(\gamma I + \tilde{C})^{-1}(\gamma I - \tilde{S})(\gamma I + \tilde{S})^{-1}\right) \\ &\leq \left\| (\gamma I - \tilde{C})(\gamma I + \tilde{C})^{-1}(\gamma I - \tilde{S})(\gamma I + \tilde{S})^{-1} \right\|_2 \\ &\leq \left\| (\gamma I - \tilde{C})(\gamma I + \tilde{C})^{-1} \right\|_2 \cdot \left\| (\gamma I - \tilde{S})(\gamma I + \tilde{S})^{-1} \right\|_2. \end{aligned} \quad (4.11)$$

Since  $C_A, C_B, F_A$  and  $F_B$  are diagonalizable by the Fourier-type matrices  $F_A, F_B, \hat{F}_A$  and  $\hat{F}_B$ , respectively, (see (2.3)), then  $\tilde{C}$  is diagonalizable by  $F_B \otimes F_A$  and  $\tilde{S}$  by  $\hat{F}_B \otimes \hat{F}_A$  [23],

$$\tilde{C} = (F_B \otimes F_A)^* \Lambda_{\tilde{C}} (F_B \otimes F_A), \quad \tilde{S} = (\hat{F}_B \otimes \hat{F}_A)^* \Sigma_{\tilde{S}} (\hat{F}_B \otimes \hat{F}_A), \quad (4.12)$$

where  $\Lambda_{\tilde{C}} = I \otimes \Lambda_A + \Lambda_B \otimes I$  and  $\Sigma_{\tilde{S}} = I \otimes \Sigma_A + \Sigma_B \otimes I$ .

Thus, by (4.12) and the invariance of the matrix 2-norm under a unitary similarity,

$$\begin{aligned} \left\| (\gamma I - \tilde{C})(\gamma I + \tilde{C})^{-1} \right\|_2 &= \left\| (\gamma I - \Lambda_{\tilde{C}})(\gamma I + \Lambda_{\tilde{C}})^{-1} \right\|_2 = \max_{\lambda_k \in \lambda(\tilde{C})} \left| \frac{\gamma - \lambda_k}{\gamma + \lambda_k} \right|, \\ \left\| (\gamma I - \tilde{S})(\gamma I + \tilde{S})^{-1} \right\|_2 &= \left\| (\gamma I - \Sigma_{\tilde{S}})(\gamma I + \Sigma_{\tilde{S}})^{-1} \right\|_2 = \max_{\mu_k \in \lambda(\tilde{S})} \left| \frac{\gamma - \mu_k}{\gamma + \mu_k} \right|. \end{aligned}$$

For  $\lambda_k = a_k + \mathbf{i}b_k \in \lambda(\tilde{C})$  and  $\mu_k = c_k + \mathbf{i}d_k \in \lambda(\tilde{S})$ ,  $k = 1, \dots, mn$ , we have  $a_k > 0$ ,  $c_k \geq 0$  (or  $a_k \geq 0$ ,  $c_k > 0$ ), by the assumption that  $\tilde{C}$  is positive definite and  $\tilde{S}$  is positive semi-definite (or vice-versa), and then

$$\left| \frac{\gamma - \lambda_k}{\gamma + \lambda_k} \right| = \sqrt{\frac{(\gamma - a_k)^2 + b_k^2}{(\gamma + a_k)^2 + b_k^2}} < 1 (\leq 1), \quad \left| \frac{\gamma - \mu_k}{\gamma + \mu_k} \right| = \sqrt{\frac{(\gamma - c_k)^2 + d_k^2}{(\gamma + c_k)^2 + d_k^2}} \leq 1 (< 1),$$

since  $\gamma > 0$ . As a consequence,

$$\sigma_\gamma := \max_{\lambda_k \in \lambda(\tilde{C})} \left| \frac{\gamma - \lambda_k}{\gamma + \lambda_k} \right| \cdot \max_{\mu_k \in \lambda(\tilde{S})} \left| \frac{\gamma - \mu_k}{\gamma + \mu_k} \right| < 1. \quad (4.13)$$

Finally, (4.11) yields

$$\rho(\mathcal{M}_\gamma) \leq \sigma_\gamma < 1$$

which ensures that the CSCS iteration (4.3) converges to the exact solution  $\mathbf{x}^*$  of the linear system  $\mathcal{A}\mathbf{x} = \mathbf{c}$  and that the CSCS iteration (3.5), which is equivalent to (4.3), converges to the exact solution  $X^* \in \mathbb{C}^{m \times n}$  of the Sylvester equation (1.1).  $\square$

**COROLLARY 4.2.** *If one of the matrices  $C_A$ ,  $C_B$ ,  $S_A$  and  $S_B$  is positive definite and all the others are positive semi-definite, then the CSCS iteration (3.5) converges to the exact solution  $X^* \in \mathbb{C}^{m \times n}$  of the Sylvester equation (1.1).*

*Proof.* Recall that given two matrices  $G \in \mathbb{C}^{n \times n}$  and  $H \in \mathbb{C}^{m \times m}$  with eigenvalues  $\lambda_i$ ,  $i = 1, \dots, n$ , and  $\mu_j$ ,  $j = 1, \dots, m$ , respectively, the eigenvalues of the Kronecker sum  $G \oplus H = I_m \otimes G + H \otimes I_n$  are the pairwise sums  $\lambda_i + \mu_j$ ,  $i = 1, \dots, n$ ,  $j = 1, \dots, m$  (see, e.g., [23]). Using this property, it is immediate to conclude, for example, that if  $G$  is positive definite and  $H$  is positive semi-definite, then  $G \oplus H$  is positive definite. In fact,  $\text{Re}(\lambda_i) > 0$  and  $\text{Re}(\mu_j) \geq 0$  imply that  $\text{Re}(\lambda_i + \mu_j) = \text{Re}(\lambda_i) + \text{Re}(\mu_j) > 0$ .

Suppose that  $C_A$  is positive definite and  $C_B$ ,  $S_A$  and  $S_B$  are all positive semi-definite. Then  $\tilde{C} = C_A \oplus C_B^T$  and  $\tilde{S} = S_A \oplus S_B^T$ , defined in (4.7), are positive definite and positive semi-definite, respectively. According to Theorem 4.1, the spectral radius  $\rho(\mathcal{M}_\gamma)$  of the iteration matrix  $\mathcal{M}_\gamma$  in (4.8) is less than 1 and the CSCS iteration (3.5) converges. All the other cases are similar.  $\square$

Next theorem addresses the issue of how to obtain a value for  $\gamma = \alpha + \beta$  (and naturally for  $\alpha$  and  $\beta$ ) that leads to a good convergence speed.

**REMARK 4.1.** Let  $\lambda_k = a_k + \mathbf{i}b_k \in \lambda(\tilde{C})$  and  $\mu_k = c_k + \mathbf{i}d_k \in \lambda(\tilde{S})$ ,  $k = 1, \dots, mn$ , satisfy

$$\theta_{min} \leq a_k, c_k \leq \theta_{max} \quad \text{and} \quad \eta_{min} \leq |b_k|, |d_k| \leq \eta_{max}, \quad (4.14)$$

where  $\theta_{min}$  and  $\theta_{max}$  are the lower and upper bounds, respectively, of the real part of the eigenvalues  $\lambda(\tilde{C}) \cup \lambda(\tilde{S})$ , and  $\eta_{min}$  and  $\eta_{max}$  are the lower and the upper bounds, respectively, of the absolute values of the imaginary part of the eigenvalues  $\lambda(\tilde{C}) \cup \lambda(\tilde{S})$ . A bound for  $\sigma_\gamma$  is given by

$$\max_{(\theta, \eta) \in \Omega} \frac{(\gamma - \theta)^2 + \eta^2}{(\gamma + \theta)^2 + \eta^2} \quad (4.15)$$

where  $\Omega = [\theta_{min}, \theta_{max}] \times [\eta_{min}, \eta_{max}]$ , since

$$\sigma_\gamma \leq \max_{\lambda_k \in \lambda(\tilde{C}) \cup \lambda(\tilde{S})} \left| \frac{\gamma - \lambda_k}{\gamma + \lambda_k} \right| \cdot \max_{\mu_k \in \lambda(\tilde{C}) \cup \lambda(\tilde{S})} \left| \frac{\gamma - \mu_k}{\gamma + \mu_k} \right| = \max_{\lambda_k \in \lambda(\tilde{S}) \cup \lambda(\tilde{C})} \left| \frac{\gamma - \lambda_k}{\gamma + \lambda_k} \right|^2.$$

We may consider that the optimal choice  $\gamma^*$  for the shift parameter  $\gamma$  is the value that minimizes the above estimate (4.15). The following theorem gives an explicit formula for  $\gamma^*$ , if  $\theta_{min} > 0$ .

**THEOREM 4.3.** *If  $\theta_{min} \geq 0$ , the minimum value*

$$\min_{\gamma > 0} \left\{ \max_{(\theta, \eta) \in \Omega} \frac{(\gamma - \theta)^2 + \eta^2}{(\gamma + \theta)^2 + \eta^2} \right\}$$

is attained at

$$\gamma^* = \begin{cases} \sqrt{\theta_{min}\theta_{max} - \eta_{max}^2} & \text{for } \eta_{max} < \tilde{\eta} \\ \sqrt{\theta_{min}^2 + \eta_{max}^2} & \text{for } \eta_{max} \geq \tilde{\eta}, \end{cases} \quad (4.16)$$

and it is equal to

$$\sigma^* = \begin{cases} \frac{\theta_{min} + \theta_{max} - 2\sqrt{\theta_{min}\theta_{max} - \eta_{max}^2}}{\theta_{min} + \theta_{max} + 2\sqrt{\theta_{min}\theta_{max} - \eta_{max}^2}} & \text{for } \eta_{max} < \tilde{\eta}, \\ \frac{\sqrt{\theta_{min}^2 + \eta_{max}^2} - \theta_{min}}{\sqrt{\theta_{min}^2 + \eta_{max}^2} + \theta_{min}} & \text{for } \eta_{max} \geq \tilde{\eta}. \end{cases}$$

where  $\tilde{\eta} = \sqrt{\theta_{min}(\theta_{max} - \theta_{min})}/2$ .

The proof of this theorem can be found in [4, pp. 324–326] and [3].

Concerning the choice of the shift parameters  $\alpha$  and  $\beta$  in the CSCS method, to choose  $\alpha = \beta = \gamma^*/2$ , where  $\gamma^*$  is computed using (4.16), seems to be a natural choice in the case that  $A$  and  $B$  have approximate norms. In practice Theorem 4.3 gives an efficient procedure to compute  $\alpha$  and  $\beta$  since we have the explicit formulae for the eigenvalues of the matrices  $C_A$ ,  $C_B$ ,  $S_A$  and  $S_B$  (see (2.4)) and we use these formulae to implement CSCS. Thus we can obtain the eigenvalues of  $\tilde{C}$  and  $\tilde{S}$  as a byproduct and verify if the sufficient condition for convergence given by Theorem 4.1 is satisfied. See Algorithm 4 in Appendix A. Notice that the case  $\theta_{min} = 0$  brings no difficulty in computing  $\gamma^*$  - when  $\theta_{min} = 0$ , we have  $\tilde{\eta} = 0$  and  $\gamma^* = \eta_{max}$ .

**5. Numerical results.** In this section we illustrate the performance of the CSCS algorithm exhibiting some numerical examples. We compare the computational behavior of this method with the Hermitian and skew-Hermitian splitting iteration (HSS) [2] and with a block variant of the Symmetric Successive Over-Relaxation scheme (BSSOR) [26, 35, 43, 42].

All the algorithms were implemented in MATLAB (R2020b) in double precision (unit roundoff  $\varepsilon = 2.2 \cdot 10^{-16}$ ) on a LAPTOP-KVSVAUU8 with an Intel(R) Core(TM) i5-8250U CPU @ 1.60GHz and 8 GB RAM, under Windows 10 Home. See Appendix A for details on the MATLAB implementations (Algorithms 1, 5 and 6 for CSCS, HSS and BSSOR, respectively). No parallel MATLAB operations were used.

The built-in functions `fft` and `ifft` (Discrete Fourier transform and its inverse) were used in CSCS, in particular to compute the residual  $R = C - AX - XB$  (see Algorithm 3 in Appendix A). The use of sparse techniques is an alternative way to compute the residual  $C - AX - XB$ . Indeed, if our matrices are stored in sparse format (even if only  $A$ ,  $B$  and  $C$ ), then MATLAB will automatically use highly efficient multiplication. The advantage of using Discrete Fourier transforms over these sparse techniques can only be observed for dense Toeplitz matrices (matrices with a low sparsity pattern or full matrices). See Example 5.4.

Hermitian and skew-Hermitian matrices can be diagonalizable by unitary matrices and thus it is possible to treat the two steps at each iteration of the HSS method very efficiently - the linear systems are all diagonal. The diagonalization process is carried out by MATLAB functions `schur` and `rsf2csf` for the real and complex Schur decompositions.

We use a variant of the BSOR (block SOR) which combines two BSOR steps together in one iteration. Specifically, BSSOR is a forward BSOR step followed by a backward BSOR step. The roles of the triangular factors  $L$  and  $U$  of both  $A$  and  $B$  are reversed in the second step. The value of the relaxation parameter  $\omega$  is the same in both steps. We remark here that the application of SSOR (Symmetric SOR) as a preconditioner for other iterative schemes, in the case of symmetric matrices, was the primary motivation for SSOR, since the convergence rate is usually slightly slower than the convergence rate of SOR with optimal  $\omega$ . In our comparison study, in particular of the number of iterations needed for convergence, it seems more appropriate to use BSSOR than BSOR given that each iteration of BSSOR consists of two steps, like CSCS and HSS.

The occurring linear systems in BSSOR are solved with the MATLAB function `linsolve` which uses  $LU$  factorization with partial pivoting when the coefficient matrix is square. This function is more efficient than the backslash operator since it is possible to specify the appropriate solver as determined by the properties of the matrix.

We also compare our method with the Bartels–Stewart direct method as implemented in the MATLAB function `lyap` from the Control Toolbox. This function performs the real Schur decompositions of  $A$  and  $B$  in equation (1.1), lower and upper, respectively, and converts them afterwards to their complex forms; computes the solution of the resulting Sylvester equation solving  $m$  triangular systems and then transforms this solution back to the solution of the original Sylvester equation. See Algorithm 7 in Appendix A for our own implementation of this method (`mylyap` function).

The null matrix was chosen as the initial approximation,  $X^{(0)} = O$ , in all our numerical experiments, and the stopping criterion implemented was

$$\frac{\|R^{(k)}\|_F}{\|C\|_F} \leq tol, \quad (5.1)$$

where  $R^{(k)} = C - AX^{(k)} - X^{(k)}B$  is the residual attained at iteration  $k$  and  $tol$  is the desired accuracy, usually set to  $10^{-6}$ .

In our first example we analyze a standard Sylvester equation that comes from a finite difference discretization of the two dimensional convection-diffusion equation

$$-(u_{xx} + u_{yy}) + \sigma(x, y)u_x + \tau(x, y)u_y = f(x, y), \quad (5.2)$$

posed on the unit square  $(0, 1) \times (0, 1)$  with Dirichlet-type boundary conditions. Here we consider the case when the coefficients  $\sigma$  and  $\tau$ , which represent the velocity components along the  $x$  and  $y$  directions, respectively, are constant. See [14, p. 371]. A five-point discretization of the operator leads to a linear system

$$\mathcal{A}\mathbf{u} = \mathbf{v}, \quad (5.3)$$

where now  $\mathbf{u}$  denotes a vector in a finite-dimensional space. We consider a uniform  $n \times n$  grid and use standard second-order finite differences for the Laplacian  $u_{xx} + u_{yy}$  and either centered or upwind differences for the first derivatives  $u_x$  and  $u_y$ . See [18, p. 217]. With  $\mathbf{u}$  ordered lexicographically in the natural ordering as  $(u_{11}, u_{21}, \dots, u_{nn})^T$ , the coefficient matrix  $\mathcal{A}$  is a block tridiagonal matrix whose  $j$ th row contains the subdiagonal, diagonal and superdiagonal blocks, all of order  $n$ , respectively,

$$\mathcal{A}_{j,j-1} = bI_n, \quad \mathcal{A}_{j,j} = \text{tridiag}(c, a, d), \quad \mathcal{A}_{j,j+1} = eI_n, \quad (5.4)$$

where  $a, b, c, d$  and  $e$  depend on the discretization. Blocks  $\mathcal{A}_{1,0}$  and  $\mathcal{A}_{n,n+1}$  are not defined. Let  $h = \frac{1}{n+1}$  ( $n$  inner grid points in each direction). After scaling by  $h^2$ , the matrix entries are given by

$$a = 4, \quad b = -\left(1 + \frac{\tau h}{2}\right), \quad c = -\left(1 + \frac{\sigma h}{2}\right), \quad d = -\left(1 - \frac{\sigma h}{2}\right), \quad e = -\left(1 - \frac{\tau h}{2}\right), \quad (5.5)$$

for the centered difference scheme, and by

$$a = 4 + (\tau + \sigma)h, \quad b = -(1 + \tau h), \quad c = -(1 + \sigma h), \quad d = -1, \quad e = -1 \quad (5.6)$$

for the upwind scheme when  $\sigma \geq 0$  and  $\tau \geq 0$ . At the  $(i, j)$  grid point, the right-hand side satisfies  $v_{ij} = h^2 f_{ij}$ , where  $f_{ij} = f(ih, jh)$ .

When  $e = d$  and  $b = c$ , the coefficient matrix  $\mathcal{A}$  in the linear system (5.3) can be written in the form  $\mathcal{A} = I_n \otimes A + A \otimes I_n$  where  $A = \text{tridiag}(c, a/2, d)$ . Therefore, the Sylvester equation

$$AX + XA^T = V \quad (5.7)$$

is equivalent to the linear system (5.3), where  $X$  and  $V$  are the matrix-stacking of the vectors  $\mathbf{u}$  and  $\mathbf{v}$ , respectively.

Different discretization schemes of equation (5.2) will naturally lead to different Sylvester equations (and different discretization errors). In [35] it is described how we can obtain a general equation  $AX + XB = C$  for any values of  $\sigma$  and  $\tau$  applying the central differences operator. Matrix  $A$  corresponds to the discretization in the  $y$ -direction and matrix  $B$  in the  $x$ -direction. When  $\sigma$  and  $\tau$  are constant,  $A$  and  $B$  are tridiagonal Toeplitz matrices defined by

$$A = \text{tridiag}\left(-1 + \frac{\tau h}{2}, 2, -1 - \frac{\tau h}{2}\right) \quad \text{and} \quad B = \text{tridiag}\left(-1 + \frac{\sigma h}{2}, 2, -1 - \frac{\sigma h}{2}\right) \quad (5.8)$$

with  $A = B$  if  $\tau = \sigma$ .

#### EXAMPLE 5.1.

Here we solve the Sylvester equation (5.7) representing the convection-diffusion equation (5.2) with homogeneous Dirichlet boundary conditions and the function  $f$  defined by  $f(x, y) = e^{x+y}$ . Different values for  $\tau = \sigma$  and the step size  $h = \frac{1}{n+1}$  are considered.

The performance of all the methods, BSSOR, HSS and CSCS, concerning the number of iterations (iter) and CPU time in seconds ( $t_{\text{CPU}}$ ) are shown in Tables 5.1, for the centered differences scheme (5.5). The results for the upwind scheme (5.6) and for the alternative scheme (5.8) are practically the same.

In the CSCS method we took  $\alpha = \beta \approx \gamma^*$ , where  $\gamma^*$  is computed using the expression (4.16), and for the HSS method we chose  $\alpha = \beta \approx 2\gamma^*$  (as a result of a numerical search around  $\gamma^*/2$ ); for the relaxation parameter  $\omega$  in the BSSOR method we used the heuristic estimate given by  $\omega = 2 - 10h$  (approximately).

We report that the initial matrix  $C_A$  is positive semi-definite but  $S_A$  is positive definite (as well as  $\tilde{C} = I_n \otimes C_A + C_A \otimes I_n$  and  $\tilde{S} = I_n \otimes S_A + S_A \otimes I_n$ , respectively), and thus the CSCS method always converges. In fact, we can prove that this splitting property of the matrix  $A = \text{tridiag}(c, a/2, d)$  is true in general for any positive values of  $\sigma$  and  $\tau$ .

For  $h = 0.05$  the BSSOR method converged but very slowly. It took more than 20 minutes to deliver a solution, with relative residual norm of about  $10^{-4}$  (1500 iterations), for  $\sigma = 2$ , and  $10^{-6}$  (946 iterations), for  $\sigma = 10$ . It is not a suitable method for this case.

Overall, the number of iterations needed for convergence by all the methods is relatively high and this reflects the fact that the spectral radii of the iteration matrices are closer to 1 than to 0. Nevertheless, CSCS exhibits the best behavior among the three methods. Our method is nearly 3 times faster than HSS (5 times for  $\sigma = 2$ ,  $n = 399$ ) and 8 times faster (in average) than BSSOR, for  $n \leq 199$  (much faster for  $n > 199$ ).

$h = \frac{1}{n+1}$	BSSOR			HSS			CSCS			
	$\omega$	iter	t <sub>CPU</sub>	$\alpha = \beta$	iter	t <sub>CPU</sub>	$\alpha = \beta$	iter	t <sub>CPU</sub>	
$\sigma = 2$	0.04	1.75	79	0.04	0.20	85	0.01	0.10	42	0.005
	0.02	1.85	167	0.25	0.10	167	0.07	0.045	84	0.03
	0.01	1.95	309	2.25	0.050	328	0.62	0.023	168	0.25
	0.005	1.95	767	51.6	0.025	648	5.41	0.011	342	1.90
	0.0025	-	-	-	0.013	1285	90.3	0.006	700	18.2
$\sigma = 10$	0.04	1.75	36	0.02	0.45	64	0.01	0.20	29	0.006
	0.02	1.85	69	0.11	0.22	126	0.06	0.075	56	0.02
	0.01	1.85	190	1.44	0.11	252	0.44	0.038	108	0.17
	0.005	1.95	258	22.6	0.05	448	3.38	0.019	216	1.22
	0.0025	-	-	-	0.013	841	66.0	0.0094	438	20.9

Table 5.1 BSSOR, HSS and CSCS performance for the Example 5.1 (centered difference scheme).

In [9, 11, 25] the authors study the numerical solution of (5.2) with non-constant coefficients. The discretization matrices  $\mathcal{A}_1$  and  $\mathcal{A}_2$  from two different linear systems (5.3) are used to create a Sylvester equation

$$\mathcal{A}_1 X + X \mathcal{A}_2 = \mathcal{C}, \quad (5.9)$$

where  $\mathcal{C}$  is randomly generated from values uniformly distributed in  $[0, 1]$ . These numerical examples were devised entirely for testing purposes and they are not connected to the solution of (5.3). We will imitate this type of examples but in our case matrices  $\mathcal{A}_1$  and  $\mathcal{A}_2$  must be Toeplitz.

EXAMPLE 5.2.

We slightly change matrix  $\mathcal{A}$ , defined by (5.4), to have constant diagonal, subdiagonal, superdiagonal,  $n^{\text{th}}$  diagonal and  $-n^{\text{th}}$  diagonal (values  $a$ ,  $c$ ,  $d$ ,  $e$  and  $b$ , respectively). For different values of  $\sigma = \tau$ , we define  $\mathcal{A}_1$  and  $\mathcal{A}_2$ , with orders  $n^2$  and  $m^2$ , respectively, and solve (5.9). Table 5.2 shows the outcome of this experiment.

$n^2; m^2$	BSSOR			HSS			CSCS			
	$\omega$	iter	t <sub>CPU</sub>	$\alpha = \beta$	iter	t <sub>CPU</sub>	$\alpha = \beta$	iter	t <sub>CPU</sub>	
$\sigma_1 = \sigma_2 = 2$	49; 100	1.75	38	0.12	0.89	49	0.05	0.60	30	0.02
	100; 100	1.75	39	0.28	0.81	65	0.13	0.41	33	0.05
	225; 225	1.75	63	4.84	0.45	92	1.20	0.27	46	0.25
	225; 400	1.85	67	24.5	0.42	99	3.92	0.28	58	0.76
	400; 400	1.85	74	61.7	0.35	117	9.20	0.20	61	1.18
	625; 625	1.75	98	333.7	0.29	143	39.1	0.19	89	4.63
$\sigma_1 = 1; \sigma_2 = 10$	100; 100	1.75	39	0.28	0.87	69	0.15	0.87	69	0.13
	225; 400	1.85	67	27.89	0.45	100	4.03	0.31	64	0.88
	625; 625	1.85	98	335.0	0.27	150	40.7	0.14	74	4.00
	784; 784	1.85	116	700.3	0.29	164	89.5	0.15	84	7.00

Table 5.2 Performance of BSSOR, HSS and CSCS for the Example 5.2.

We verified that the matrices  $\tilde{C}$  and  $\tilde{S}$  are positive semi-definite and positive definite, respectively, like in the first example. The values for the parameters  $\alpha$  and  $\beta$  that led to a smaller number of iterations were values greater than the value  $\gamma^*$  given by (4.16), by a factor of about 10 or higher.

The convergence rate of all the methods is faster for this example than for the previous one and, as expected, BSSOR is a very slow method compared to HSS and CSCS. Also in this case, when compared to HSS, the CSCS method is about 5 times faster, for matrices of order 200, and 8 times faster if the order of the matrices raises above 600.

The next numerical example can be found, for instance, in [2, 41, 44]. As mentioned in [2] this class of problems appears associated with the preconditioned Krylov subspace iteration method used to solve the systems of linear equations which arise from the discretization of various differential equations and boundary value problems using finite difference or Sinc-Galerkin schemes.

EXAMPLE 5.3. Consider the Sylvester equation (1.1) with matrices  $A, B \in \mathbb{C}^{n \times n}$  ( $m = n$ ) defined by

$$A = B = M + 2rN + \frac{100}{(n + 1)^2}I,$$

where  $M, N \in \mathbb{C}^{n \times n}$  are Toeplitz tridiagonal matrices,  $M = \text{tridiag}(-1, 2, -1)$ ,  $N = \text{tridiag}(0.5, 0, -0.5)$ . In a more compacted form,

$$A = B = \text{tridiag}\left(-1 + r, 2 + \frac{100}{(n + 1)^2}, -1 - r\right).$$

The parameter  $r$  depends on the properties of the problem being discretized.

Although this problem is similar to the one considered in Example 5.1, we decided to show the results of our experiments in order to compare them with the results presented by other authors, namely in [2, 29, 41]. Table 5.3 contains the summary of our experiments for different instances of the parameter  $r$  and the order  $n$  of the matrices.

	$n$	BSSOR			HSS			CSCS		
		$\omega$	iter	$t_{\text{CPU}}$	$\alpha$	iter	$t_{\text{CPU}}$	$\alpha = \beta$	iter	$t_{\text{CPU}}$
$r = 0.01$	64	1.75	36	0.10	0.17	123	0.10	0.130	32	0.02
	128	1.85	71	1.08	0.09	244	0.75	0.070	60	0.14
	256	1.95	167	24.0	0.05	453	19.2	0.035	112	0.88
	512	1.95	281	670.3	0.05	520	64.0	0.017	221	7.00
	1024	-	-	-	0.01	1204	1408	0.010	392	61.0
$r = 0.1$	64	1.75	35	0.12	0.23	90	0.07	0.14	31	0.02
	128	1.85	65	1.10	0.13	145	0.46	0.08	55	0.13
	256	1.85	139	21.4	0.09	219	4.26	0.05	86	0.70
	512	1.75	455	678.1	0.10	314	33.4	0.10	317	10.4
	1024	-	-	-	0.10	607	636.4	0.10	610	100.5
$r = 1$	64	1.5	22	0.06	0.81	40	0.04	0.26	26	0.01
	128	1.5	33	0.50	0.62	60	0.19	0.16	41	0.09
	256	1.75	47	7.10	0.51	92	1.81	0.11	61	0.45
	512	1.75	62	114.2	0.25	132	16.7	0.25	138	3.95
	1024	-	-	-	0.25	192	173.0	0.15	171	21.3

Table 5.3 Performance of BSSOR, HSS and CSCS methods for the Example 5.3.

The values of the shift parameter  $\alpha$  in the HSS method are the values which were presented in [2], for  $n \leq 256$  (see  $\omega_{\text{exp}}$  and  $\alpha_{\text{exp}}$  in [2, Table 4.2], obtained through an experimental search). The values given to the shift parameters  $\alpha$  and  $\beta$  in the CSCS method were determined using the expression (4.16) - we computed  $\gamma^*$  and let  $\alpha = \beta$  between  $\gamma^*/8$  and  $\gamma^*/2$  - and these values are also used with HSS when  $n > 256$ .

In this example the convergence is faster than in Example 5.1, in particular when  $r = 1$ . Matrices  $\tilde{C}$  and  $\tilde{S}$  are both positive definite and the CSCS method outperforms the HSS and BSSOR methods both in terms of the number of iterations and in what respects to the computational efficiency. BSSOR may be very slow for matrices of order  $n \geq 512$ , taking more than 20 minutes to converge. Compared to HSS the CPU time required by CSCS to converge is, in most cases, 4 to 8 times smaller (in extreme cases, this factor may be much smaller). Except for  $r = 1$ , our implementation of HSS demands a higher number of iterations than shown in [2] for this same method, but despite this, in all cases the CPU time needed is reduced.

The advantage of using FFT operations in the CSCS method can be entirely appreciated when we take  $A$  and  $B$  to be full Toeplitz matrices. Next example considers this case and reports the CPU elapsed times for CSCS and MATLAB function `lyap`.

EXAMPLE 5.4. *This example takes positive definite circulant and skew-circulant matrices  $C_A$  and  $S_A$  (obtained using translation of origin on randomly generated matrices) and forms  $A = C_A + S_A$ ,  $B = A$ . Matrix  $C$  is chosen to be the matrix attained when all the entries in  $X$  are set to be 1.*

*We take  $\alpha = \beta = \gamma^*/2$  where  $\gamma^*$  is computed using the expression (4.16) in Theorem 4.3. See Table 5.4 for a comparison of the efficiency of CSCS and `lyap`.*

$n$	CSCS			<code>lyap</code>		
	$\alpha = \beta$	iter	resid	$t_{\text{CPU}}$	resid	$t_{\text{CPU}}$
100	43.49	5	$1.1 \cdot 10^{-6}$	0.017	$2.2 \cdot 10^{-15}$	0.013
		12	$1.9 \cdot 10^{-15}$	0.040		
250	103.75	5	$2.04 \cdot 10^{-6}$	0.07	$1.8 \cdot 10^{-15}$	0.10
		13	$1.5 \cdot 10^{-15}$	0.14		
500	208.0	5	$2.0 \cdot 10^{-6}$	0.32	$1.9 \cdot 10^{-15}$	0.26
		12	$4.0 \cdot 10^{-15}$	0.67		
1000	426.6	5	$1.3 \cdot 10^{-6}$	1.42	$2.1 \cdot 10^{-15}$	1.27
		12	$9.4 \cdot 10^{-15}$	3.06		
1500	645.4	5	$1.2 \cdot 10^{-6}$	3.20	$2.3 \cdot 10^{-15}$	3.34
		13	$3.7 \cdot 10^{-16}$	7.43		
2000	856.97	5	$1.7 \cdot 10^{-6}$	5.62	$2.5 \cdot 10^{-15}$	7.01
		12	$8.2 \cdot 10^{-15}$	12.01		
2500	1080.1	5	$1.7 \cdot 10^{-6}$	10.64	$2.7 \cdot 10^{-15}$	16.01
		13	$1.0 \cdot 10^{-15}$	23.78		

Table 5.4 Performance of CSCS and `lyap` for full matrices  $A$  and  $B$ .

If the relative accuracy demanded is  $\mathcal{O}(10^{-6})$ , which is often enough in many applications, the CSCS method is comparable to or even faster than `lyap`. When full accuracy  $\mathcal{O}(\varepsilon)$  is important, more iterations are needed and CSCS takes approximately twice as long as `lyap`, which, however, can still be considered very satisfactory since these methods are fast even for large dimensions like  $n \geq 1000$ .

We may take our function `mylyap` (see Algorithm 7) in this comparison study, which is possibly the fairest comparison study to present, given that in our implementations we are not capable of reproducing the MATLAB internal linear systems solvers used by `lyap`. We clearly acknowledge that, when full accuracy  $\mathcal{O}(\varepsilon)$  is required, CSCS method is always faster, about 10 times faster, than `mylyap` for full Toeplitz matrices

$A$  and  $B$  (`mylyap` is, as expected, slower than `lyap`).

**6. Conclusions.** We considered the problem of solving a large continuous Sylvester equation  $AX + XB = C$  where the coefficient matrices  $A$  and  $B$  are assumed to be Toeplitz matrices and we have devised the CSCS iteration which is a method based on the circulant and skew-circulant splittings of the matrices  $A$  and  $B$ . The spectral properties of these structured matrices allow the use of fast Fourier transforms (FFTs) which reduces significantly the operation count of matrix multiplication and thus the computational efficiency of the algorithm. We have also analyzed sufficient conditions for the convergence of the CSCS iteration and have derived an upper bound for its convergence factor. The numerical experiments we have carried out illustrate that CSCS is a faster and more robust iterative algorithm than the alternatives HSS and BSSOR. The advantage of using FFT operations in the CSCS method can be entirely appreciated when we take  $A$  and  $B$  to be full Toeplitz matrices and in this case CSCS is a very competitive algorithm even when compared with the MATLAB function `lyap` which implements the Bartels–Stewart direct method. Moreover, since FFT-based operations have very high parallel potentialities, our CSCS algorithm is therefore suited for parallel frameworks.

## REFERENCES

- [1] Z.-Z. Bai, *Splitting iteration methods for non-Hermitian positive definite systems of linear equations*, Hokkaido Math. J., 36 (2007), pp. 801-814.
- [2] Z.-Z. Bai, *On Hermitian and skew-Hermitian splitting iterative methods for continuous Sylvester equations*, J. Comput. Math., 29 (2011), pp. 185-198.
- [3] Z.-Z. Bai and M.-K. Ng, *Erratum*, Numer. Linear Algebra Appl., 19 (2012), p. 891.
- [4] Z.-Z. Bai, G. H. Golub and M.-K. Ng, *On successive overrelaxation acceleration of the Hermitian and skew-Hermitian splitting iterations*, Numer. Linear Algebra Appl. 14 (2007), pp. 319-335.
- [5] Z.-Z. Bai, G. H. Golub and M. K. Ng, *Hermitian and skew-Hermitian splitting methods for non-Hermitian positive definite linear systems*, SIAM J. Matrix Anal. Appl., 24 (2003), pp. 603-626.
- [6] Z.-Z. Bai, X.-X. Guo and S.-F. Xu, *Alternately linearized implicit iteration methods for the minimal nonnegative solutions of the nonsymmetric algebraic Riccati equations*, Numer. Linear Algebra Appl., 13(2006), pp. 655-674.
- [7] R. H. Bartels and G. W. Stewart, *Solution of the matrix equation  $AX + XB = C$ : Algorithm 432*, Commun. ACM, 15 (1972), pp. 820-826
- [8] P. Benner, R. C. Li and N. Truhar, *On the ADI method for Sylvester equations*, J. Comput. Appl. Math., 233 (2009), pp. 1035-1045.
- [9] P. Benner and P. Kürschner, *Computing real low-rank solutions of Sylvester equations by the factored ADI method*, Comput. Math. with Appl., 67 (9) (2014), pp. 1656-1672.
- [10] R. Bhatia and P. Rosenthal, *How and why to solve the operator equation  $AX+XB = Y$* , Bull. Lond. Math. Soc., 29 (1997), pp. 1-21.
- [11] A. Bouhamidi, M. Hached, M. Heyouni and K. Jbilou, *A preconditioned block Arnoldi method for large Sylvester matrix equations*, Numer. Linear Algebra Appl., 20 (2013), pp. 208-219.
- [12] D. Calvetti and L. Reichel, *Application of ADI iterative methods to the restoration of noisy images*, SIAM J. Matrix Anal. Appl., 17 (1996), pp. 165-186.
- [13] R. Chan and M. Ng, *Conjugate gradient methods for Toeplitz systems*, SIAM Rev., 38 (1996), pp. 427-482.
- [14] Y. H. Chen and T. W. H. Sheu, *Two-dimensional scheme for convection-difusion with linear production*, Numerical Heat Transfer, Part B, 37 (2000), pp. 365-377.
- [15] P. J. Davis, *Circulant Matrices*, John Wiley, New York, 1979.
- [16] P. Van Dooren, *Structured linear algebra problems in digital signal processing*, Numerical Linear Algebra, Digital Signal Processing and Parallel Algorithms, NATO Series F, Springer, (1991), pp. 361-384.
- [17] G.-R. Duan, *Generalized Sylvester Equations, Unified Parametric Solutions, 1st edition*, CRC Press, 2020.
- [18] H. C. Elman and G. H. Golub, *Iterative methods for cyclically reduced on-self-adjoint linear systems II.*, Math. Comput., 56(193), (1991), pp. 215-242
- [19] M. Epton, *Methods for the solution of  $AXD - BXC = E$  and its application in the numerical solution of implicit ordinary differential equations*, BIT, 20 (1980), pp. 341-345.
- [20] G. H. Golub and C. F. Van Loan, *Matrix Computations, 3rd edition*, Johns Hopkins University Press, Baltimore, Maryland, 1996.
- [21] G. H. Golub, S. G. Nash and C. F. Van Loan, *A Hessenberg-Schur method for the problem  $AX + XB = C$* , IEEE Trans. Automat. Control, 24 (1979), 909-913.
- [22] C.-Q. Gu and H.-Y. Xue, *A shift-splitting hierarchical identification method for solving Lyapunov matrix equations*, Linear

- Algebra Appl., 430 (2009), pp. 1517–1530.
- [23] R. A. Horn and C. R. Johnson, *Topics in Matrix Analysis*, Cambridge University Press, Cambridge, UK, 1991.
  - [24] D.-Y. Hu and L. Reichel, *Krylov-subspace methods for the Sylvester equation*, Linear Algebra Appl., 172 (1992), pp. 283–313.
  - [25] K. Jbilou, *Low rank approximate solutions to large Sylvester matrix equations*, Appl. Math. Comput., 177 (2006), pp. 365–376.
  - [26] S. Kadry and Z. Woznicki, *On Discussion of SOR Method for Solving the Sylvester Equation*, Int. J. Soft Comput., 2(2) (2007), 236–242.
  - [27] A. Kittisopaporn and P. Chansangiam, *Approximated least-squares solutions of a generalized Sylvester-transpose matrix equation via gradient-descent iterative algorithm*, Adv. Differ. Equ., 266 (2021). <https://doi.org/10.1186/s13662-021-03427-4>
  - [28] P. Lancaster and M. Tismenetsky, *The Theory of Matrices*, 2nd edition, Academic Press, Orlando, 1985.
  - [29] Z. Y. Liu, Y. Zhou and Y. L. Zhang, *On inexact ADI iteration for continuous Sylvester equations*, Numerical Linear Algebra with Applications, 27(5), e2320, (2020).
  - [30] Z. Y. Liu, F. Zhang, Y. Zhou, C. Ferreira and Y.L. Zhang, *Extrapolated and successive overrelaxation ADI methods for continuous Sylvester equations*, submitted.
  - [31] N. Levenberg and L. Reichel, *A generalized ADI iterative method*, Numer. Math., 66 (1993), pp.215-233.
  - [32] C. Van Loan, *Computational Frameworks for the Fast Fourier Transform*, Frontiers in Applied Mathematics, SIAM, Philadelphia, PA, USA (1992).
  - [33] M. K. Ng, *Circulant and skew-circulant splitting methods for Toeplitz systems*, Journal of Computational and Applied Mathematics 159 (2003), pp. 101-108
  - [34] M. K. Ng, *Iterative methods for Toeplitz systems*, Oxford Univ. Press, 2004.
  - [35] G. Starke and W. Niethammer, *SOR for  $AX - XB = C$* , Linear Algebra Appl., 154/156 (1991) 355–375.
  - [36] Y. Saad, *Iterative methods for sparse Linear systems*, 2nd edition, SIAM, Philadelphia, PA, USA, 2003.
  - [37] R. A. Smith, *Matrix equation  $XA + BX = C$* , SIAM J. Appl. Math., 16 (1968), 198–201.
  - [38] V. Simoncini, *Computational Methods for Linear Matrix Equations*, SIAM Rev., 8(3), (2016), pp. 377-441.
  - [39] M. Vetterli and H. Nussbaumer, *Simple FFT and DCT algorithms with reduced number of operations*, Signal Process., 6 (1984), pp. 267-278.
  - [40] E. L. Wachspress, *Iterative solution of the Lyapunov matrix equation*, Appl. Math. Lett., 107(1), (1988), pp. 87–90.
  - [41] X. Wang, W.-W. Li and L.-Z. Mao, *On positive-definite and skew-Hermitian splitting iterative methods for continuous Sylvester equation  $AX + XB = C$* , Comput. Math. Appl., 66 (2013), pp. 2352–2361.
  - [42] D. M. Young, Jr., *Iterative Solution of Large Linear Systems*, Academic Press, 1971 (reprinted by Dover, 2003).
  - [43] D. M. Young, *Convergence Properties of the Symmetric and Unsymmetric Successive Overrelaxation Methods and Related Methods*, Math. Comput., 24(112) (1970), pp. 793-807
  - [44] Q.-Q. Zheng, C.-F. Ma, *On normal and skew-Hermitian splitting iterative methods for large sparse continuous Sylvester equations*, J. Comput. Appl. Math., 268 (2014), pp. 145-154.

## Appendix A. Implementation details.

**Algorithm 1** CSCS – circulant and skew-circulant splitting iteration

**Input:** Toeplitz matrices  $A, B, C$  (orders  $n \times n, m \times m$  and  $n \times m$ ),  
initial approximation  $X_0$ , relative residual tolerance  $tol$ , maximum number of iterations  $maxit$   
**Output:** Solution  $X$  of the Sylvester equation  $AX + XB = C$

---

```

[CA, SA] = CSsplitting(Acol1, Arow1)           ▷ circulant and skew-circulant splittings of A and B
[CB, SB] = CSsplitting(Bcol1, Brow1)

Dn = exp((0 : n - 1)/n * pi * i)              ▷ Dn = [1, eπi/n, ..., e(n-1)πi/n]
Dm = exp((0 : m - 1)/m * pi * i)
DcA = diag(iff(fft(CA).').')                 ▷ DcA = diag(ΛA); ΛA = FA CA FA*
DcB = diag(iff(fft(CB).').')                 ▷ DcB = diag(ΛB); ΛB = FB CB FB*
Dsa = diag(fft(iff(Dn.' * SA * Dn).').')     ▷ Dsa = diag(ΣA); ΣA = FA SA FA*
Dsb = diag(fft(iff(Dm.' * SB * Dm).').')     ▷ Dsb = diag(ΣB); ΣB = FB SB FB*
                                           ▷ .* for element-wise product

[α, β] = shifts(DcA, DcB, Dsa, Dsb)          ▷ shift parameters α and β
D1 = DcA + α; D2 = DcB + β; D3 = Dsa + α; D4 = Dsb + β

A = sparse(A); B = sparse(B); C = sparse(C)  ▷ MATLAB sparse matrix storage format
X = X0
R = resid(X, C, DcA, DcB, Dsa, Dsb, Dn, Dm)   ▷ R = C - AX - XB using (3.12)

normR = norm(R, 'fro'); normC = norm(C, 'fro')  ▷ Frobenious norms of R and C
iter = 0                                       ▷ number of iterations counter
while ((normR/normC) > tol and iter < maxit) do
    % First step
    R = (iff(fft(R).').') * (sqrt(m/n))        ▷ transform (3.13)
    Z = R./(D1 + D2.'')                        ▷ solve (αI + ΛA)Z + Z(βI + ΛB) = R
    X = X + iff((fft(Z).').')                 ▷ update X; X ← X + FA* Z FB

    % Second step
    R = resid(X, C, DcA, DcB, Dsa, Dsb, Dn, Dm)  ▷ use (3.12) to compute R
    R = iff((fft(((Dn.') * R * Dm).').') * (sqrt(n/m)))  ▷ transform (3.14)
    Z = R./(D3 + D4.'')                       ▷ solve (αI + ΣA)Z + Z(βI + ΣB) = R
    X = X + (Dn.') * (iff(fft(Z).').') * conj(Dm)  ▷ update X; X ← X + FA* Z FB

    R = resid(X, C, DcA, DcB, Dsa, Dsb, Dn, Dm)
    normR = norm(R, 'fro')
    iter = iter + 1
end while
if iter >= maxit then
    disp('Maximum number of iterations exceed.')
end if

```

---

---

**Algorithm 2** CSsplitting – circulant and skew-circulant splitting of a Toeplitz matrix
 

---

**Input:** First column and first row,  $c$  and  $r$ , of a  $n \times n$  Toeplitz matrix  $A$

( $c$  and  $r$  should be given as rows,  $c(1)$  should be equal to  $r(1)$ )

**Output:**  $C_A$  and  $S_A$  with  $A = C_A + S_A$ , circulant and skew-circulant splitting of  $A$

```
% Circulant part
CAc = (c + [0, fliplr(r(2 : n))])/2
CAr = [r(1)/2, fliplr(CAc(2 : n))]
CA = toeplitz(CAc, CAr)
```

```
% Skew-circulant part
SAc = (c - [0, fliplr(r(2 : n))])/2
SAr = [r(1)/2 - fliplr(SAc(2 : n))]
SA = toeplitz(SAc, SAr)
```

---



---

**Algorithm 3** resid – residual for a given approximation  $X$ 


---

**Input:** Approximation  $X$  and matrix  $C$  of  $AX + XB = C$ ,

$D_{cA}$ ,  $D_{cB}$ ,  $D_{sA}$ ,  $D_{sB}$ ,  $D_n$  and  $D_m$ , computed in CSCS function

**Output:** Residual  $R = C - AX - XB$  using (3.12)

```
p1 = ifft(D1 .* fft(X))
p2 = (Dn.' .* (fft(D3 .* ifft((Dn.' .* X))))
p3 = fft(((ifft(X.').' .* (D2.')))).'
p4 = (fft((X .* Dm).')).'
p4 = (ifft((p4 .* (D4.')))).' .* (conj(Dm))
R = C - p1 - p2 - p3 - p4
```

$\triangleright p_1 = F_A^* \Lambda_A F_A X$   
 $\triangleright p_2 = \hat{F}_A^* \Sigma_A \hat{F}_A X$   
 $\triangleright p_3 = X F_B^* \Lambda_B F_B$   
 $\triangleright p_4 = X \hat{F}_B^* \Sigma_B \hat{F}_B$

---



---

**Algorithm 4** shifts – find  $\gamma^*$  and shift parameters  $\alpha$  and  $\beta$ 


---

**Input:**  $D_{cA}$ ,  $D_{cB}$ ,  $D_{sA}$ ,  $D_{sB}$  eigenvalues of  $C_A$ ,  $C_B$ ,  $S_A$ ,  $S_B$ , respectively

**Output:** shifts  $\alpha$  and  $\beta$

$D_{\tilde{C}} = D_{cA} + D_{cB}.$ ;  $D_{\tilde{S}} = D_{sA} + D_{sB}.$   $\triangleright$  eigenvalues of  $\tilde{C}$  and  $\tilde{S}$

$D_1 = [\text{real}(D_{\tilde{C}}); \text{real}(D_{\tilde{S}})];$   $D_2 = \text{abs}([\text{imag}(D_{\tilde{C}}); \text{imag}(D_{\tilde{S}})])$   $\triangleright$  see Remark 4.1

$\theta_{\min} = \min(\min(D_1));$   $\theta_{\max} = \max(\max(D_1))$

$\eta_{\min} = \min(\min(D_2));$   $\eta_{\max} = \max(\max(D_2))$

**if**  $\theta_{\min} \geq 0$  **then**  $\triangleright$  see Theorem 4.3

**if**  $\eta_{\max} < \sqrt{\theta_{\min} * (\theta_{\max} - \theta_{\min})/2}$  **then**

$\gamma^* = \sqrt{\theta_{\min} \theta_{\max} - \eta_{\max}^2}$

**else**

$\gamma^* = \sqrt{\theta_{\min}^2 + \eta_{\max}^2}$

**end if**

**else**

$\gamma^* = 1$

$\triangleright$  random value for  $\gamma^*$  if  $\theta_{\min} < 0$

$\text{disp}(\text{'Warning: } \theta_{\min} < 0. \text{ We let } \gamma^* = 1.\text{'})$

**end if**

$\alpha = \gamma^*/2;$   $\beta = \alpha$

---

**Algorithm 5 HSS – Hermitian and skew-Hermitian splitting iteration**

**Input:** Toeplitz matrices  $A, B, C$  (orders  $n \times n, m \times m$  and  $n \times m$ ),  
 initial approximation  $X_0$ , shift parameters  $\alpha$  and  $\beta$ , relative residual tolerance  $tol$ ,  
 maximum number of iterations  $maxit$

**Output:** Solution  $X$  of the Sylvester equation  $AX + XB = C$

▷ Hermitian and skew-Hermitian splittings of  $A$  and  $B$

$$H_1 = (A + A')/2$$

$$S_1 = (A - A')/2$$

$$H_2 = (B + B')/2$$

$$S_2 = (B - B')/2$$

▷ schur forms of  $H_1, H_2, S_1$  and  $S_2$  (diagonalizable)

$$[Q_1, D_1] = \text{schur}(\text{full}(H_1))$$

$$\triangleright H_1 = Q_1 D_1 Q_1^*$$

$$[Q_2, D_2] = \text{schur}(\text{full}(H_2))$$

$$\triangleright H_2 = Q_2 D_2 Q_2^*$$

$$[Q_3, D_3] = \text{schur}(\text{full}(S_1))$$

$$[Q_3, D_3] = \text{rsf2csf}(Q_3, D_3)$$

$$\triangleright S_1 = Q_3 D_3 Q_3^*$$

$$[Q_4, D_4] = \text{schur}(\text{full}(S_2))$$

$$[Q_4, D_4] = \text{rsf2csf}(Q_4, D_4)$$

$$\triangleright S_2 = Q_4 D_4 Q_4^*$$

▷ diagonal elements of  $D_1 + \alpha I_n, D_2 + \beta I_m, D_3 + \alpha I_n$  and  $D_4 + \beta I_m$

$$D_1 = \text{diag}(D_1) + \alpha; \quad D_2 = \text{diag}(D_2) + \beta$$

$$D_3 = \text{diag}(D_3) + \alpha; \quad D_4 = \text{diag}(D_4) + \beta$$

$$A = \text{sparse}(A); \quad B = \text{sparse}(B); \quad C = \text{sparse}(C)$$

▷ MATLAB sparse matrix storage format

$$X = X_0$$

$$R = C - A * X - X * B$$

$$\text{norm}R = \text{norm}(R, 'fro'); \quad \text{norm}C = \text{norm}(C, 'fro')$$

▷ Frobenious norms of  $R$  and  $C$

$$\text{iter} = 0$$

▷ number of iterations counter

**while**  $((\text{norm}R/\text{norm}C) > tol$  and  $\text{iter} < maxit)$  **do**

  % First step

$$R = Q_1' * R * Q_2$$

$$Z = R ./ (D_1 + D_2)'$$

▷ solve  $(\alpha I + D_1)Z + Z(\beta I + D_2) = R$

$$X = X + Q_1 * Z * Q_2'$$

▷ update  $X$ ;  $X \leftarrow X + Q_1 Z Q_2^*$

  % Second step

$$R = C - A * X - X * B$$

$$R = Q_3' * R * Q_4$$

$$Z = R ./ (D_3 + D_4)'$$

▷ solve  $(\alpha I + D_3)Z + Z(\beta I + D_4) = R$

$$X = X + Q_3 * Z * Q_4'$$

▷ update  $X$ ;  $X \leftarrow X + Q_3 Z Q_4^*$

$$R = C - A * X - X * B$$

$$\text{norm}R = \text{norm}(R, 'fro')$$

$$\text{iter} = \text{iter} + 1$$

**end while**

**if**  $\text{iter} \geq maxit$  **then**

$\text{disp}('Maximum number of iterations exceed.')$

**end if**

---

**Algorithm 6 BSSOR – Block Symmetric Successive Over-Relaxation iteration**


---

**Input:** matrices  $A, B, C$  (orders  $n \times n, m \times m$  and  $n \times m$ ),

initial approximation  $X_0$ , relaxation parameter  $\omega$ , relative residual tolerance  $tol$ ,

maximum number of iterations  $maxit$

**Output:** Solution  $X$  of the Sylester equation  $AX + XB = C$

`% Two steps with the same parameter  $\omega$`

`%  $(D_1/\omega + L_1)X_{k+\frac{1}{2}} + X_{k+\frac{1}{2}}(D_2/\omega + U_2) = C + [(1-w)/wD_1 - U_1]X_k + X_k[(1-w)/wD_2 - L_2]$`

`%  $(D_1/\omega + U_1)X_{k+1} + X_{k+1}(D_2/\omega + L_2) = C + [(1-w)/wD_1 - L_1]X_{k+\frac{1}{2}} + X_{k+\frac{1}{2}}[(1-w)/wD_2 - U_2]$`

`% Using residuals,  $R_k$  and  $R_{k+\frac{1}{2}}$ , and new variables,  $Z_{k+1}$  and  $Z_{k+\frac{1}{2}}$`

`%  $(D_1/\omega + L_1)Z_{k+\frac{1}{2}} + Z_{k+\frac{1}{2}}(D_2/\omega + U_2) = R_k; X_{k+\frac{1}{2}} = X_k + Z_{k+\frac{1}{2}}$`

`%  $(D_1/\omega + U_1)Z_{k+1} + Z_{k+1}(D_2/\omega + L_2) = R_{k+\frac{1}{2}}; X_{k+1} = X_{k+\frac{1}{2}} + Z_{k+1}$`

`$D_1 = \text{diag}(A); L_1 = \text{tril}(A, -1); U_1 = \text{triu}(A, 1)$`

`$\triangleright$  diagonal, strictly lower and`

`$D_2 = \text{diag}(B); L_2 = \text{tril}(B, -1); U_2 = \text{triu}(B, 1)$`

`$\triangleright$  strictly upper parts of  $A$  and  $B$`

`$D_1 = D_1/\omega; L_1 = \text{diag}(D_1) + L_1; U_1 = \text{diag}(D_1) + U_1$`

`$D_2 = D_2/\omega; L_2 = \text{diag}(D_2) + L_2; U_2 = \text{diag}(D_2) + U_2$`

`$A = \text{sparse}(A); B = \text{sparse}(B); C = \text{sparse}(C)$`

`$\triangleright$  MATLAB sparse matrix storage format`

`$X = X_0; R = C - A * X - X * B$`

`$\triangleright$  initial residual`

`$normR = \text{norm}(R, 'fro'); normC = \text{norm}(C, 'fro')$`

`$\triangleright$  Frobenious norms of  $R$  and  $C$`

`$Z = \text{zeros}(n, m)$`

`$\triangleright$  preallocation for speed`

`$\text{index}_1 = \text{eye}(n, 'logical')$`

`$\triangleright$  logical indexing for the diagonal elements`

`$\text{index}_2 = \text{eye}(m, 'logical')$`

`$\triangleright$  logical indexing for the diagonal elements`

`$\text{iter} = 0$`

`$\triangleright$  number of iterations counter`

`while  $((normR/normC) > tol$  and  $\text{iter} < maxit)$  do`

`% First step`

`$\triangleright$  solve  $L_1Z + ZU_2 = R$`

`$\text{opts.LT} = \text{true}; \text{opts.UT} = \text{false};$`

`$\triangleright$  LT - lower triangular option to linsolve`

`$L_1(\text{index}_1) = D_1 + D_2(1)$`

`$\triangleright L_1 = L_1 + D_2(1)I_n$`

`$Z(:, 1) = \text{linsolve}(L_1, R(:, 1), \text{opts})$`

`$\triangleright$  column 1 of  $Z$`

`for  $k = 2 : m$  do`

`$\triangleright$  columns 2 through  $m$  of  $Z$`

`$L_1(\text{index}_1) = D_1 + D_2(k)$`

`$\triangleright L_1 = L_1 + D_2(k)I_n$`

`$Z(:, k) = \text{linsolve}(L_1, R(:, k) - Z(:, 1 : (k-1)) * U_2(1 : (k-1), k), \text{opts})$`

`end for`

`$X = X + Z; R = C - A * X - X * B$`

`$\triangleright$  update  $X$  and  $R$`

`% Second step`

`$\triangleright$  solve  $U_1Z + ZL_2 = R$`

`$\text{opts.LT} = \text{false}; \text{opts.UT} = \text{true};$`

`$\triangleright$  UT - upper triangular option to linsolve`

`$L_2(\text{index}_2) = D_2 + D_1(n)$`

`$\triangleright L_2 = L_2 + D_1(n)I_n$`

`$Z(n, :) = \text{linsolve}(L_2.', R(n, :).', \text{opts})$`

`$\triangleright$  last row of  $Z$`

`for  $k = n - 1 : -1 : 1$  do`

`$\triangleright$  rows  $n - 1$  through 1 of  $Z$`

`$L_2(\text{index}_2) = D_2 + D_1(k)$`

`$\triangleright L_2 = L_2 + D_1(k)I_n$`

`$Z(k, :) = \text{linsolve}(L_2.', (R(k, :) - U_1(k, k+1 : n) * Z(k+1 : n, :)).', \text{opts})$`

`end for`

`$X = X + Z; R = C - A * X - X * B$`

`$\triangleright$  update  $X$  and  $R$`

`$normR = \text{norm}(R, 'fro')$`

`$\text{iter} = \text{iter} + 1$`

`end while`

`if  $\text{iter} \geq maxit$  then`

`$\text{disp}('Maximum number of iterations exceed.')$`

`end if`

---

---

**Algorithm 7** mylyap – Bartels-Stewart method
 

---

**Input:** matrices  $A, B, C$  (orders  $n \times n, m \times m$  and  $n \times m$ )

**Output:** Solution  $X$  of the Sylvester equation  $AX + XB = C$

```

[Q1, T1] = schur(full(A'))
[Q1, T1] = rsf2csf(Q1, T1)
T1 = T1'                                ▷ lower complex schur form of A; A = Q1T1Q1*

[Q2, T2] = schur(full(B))
[Q2, T2] = rsf2csf(Q2, T2)             ▷ upper complex schur form of B; B = Q2T1Q2*

dT1 = diag(T1)                          ▷ diagonal elements of T1 and T2
dT2 = diag(T2)
index = eye(n, 'logical')

% Solution of T1X + XT2 = Q1*C*Q2
C = Q1' * C * Q2
X = zeros(n, m)                          ▷ preallocation for speed
opts.LT = true                            ▷ LT - lower triangular option to linsolve
T1(index) = diag(T1) + dT2(1)             ▷ T1 = T1 + dT2(1)In
X(:, 1) = linsolve(T1, C(:, 1), opts)     ▷ column 1 of X

for k = 2 : m do                          ▷ columns 2 through m of X
    T1(index) = dT1 + dT2(k)              ▷ T1 = T1 + dT2(k)In
    X(:, k) = linsolve(T1, C(:, k) - X(:, 1 : (k - 1)) * T2(1 : (k - 1), k), opts)
end for

% Solution of AX + XB = C
X = Q1 * X * Q2'

```

---