# ON INEXACT ALTERNATING DIRECTION IMPLICIT ITERATION FOR CONTINUOUS SYLVESTER EQUATIONS

ZHONG-YUN  LIU*, YANG ZHOU*, AND YULIN ZHANG†

**Abstract.** In this paper, we study the alternating direction implicit (ADI) iteration for solving the continuous Sylvester equation $AX + XB = C$, where the coefficient matrices $A$ and $B$ are assumed to be positive semi-definite matrices (not necessarily Hermitian), and at least one of them to be positive definite. We first analyze the convergence of the ADI iteration for solving such a class of Sylvester equations, then derive an upper bound for the contraction factor of this ADI iteration. To reduce its computational complexity, we further propose an inexact variant of the ADI iteration, which employs some Krylov subspace methods as its inner iteration processes at each step of the outer ADI iteration. The convergence is also analyzed in detail. The numerical experiments are given to illustrate the effectiveness of both ADI and inexact ADI iterations.

**Key words.** Continuous Sylvester equations, ADI iteration, Inexact ADI iteration, Positive definite matrices, Convergence

**AMS subject classifications.** 15A24, 15A30, 65F10, 65F30, 65H10

**1. Introduction.** Consider the iterative solution to the continuous Sylvester equation

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

by the ADI-like iterations, where $A \in \mathbb{C}^{m \times m}$, $B \in \mathbb{C}^{n \times n}$, $C \in \mathbb{C}^{m \times n}$ are large sparse matrices.

For definiteness, throughout this paper, both $A$ and $B$ in (1.1) are assumed to be positive semi-definite† and at least one of them to be positive definite.

It is known that the Sylvester equation (1.1) has a unique solution if and only if $A$ and $-B$ haven't the common eigenvalues, see e.g., [19, 21]. A Lyapunov equation is a special case of the Sylvester equation with $B = A^H$ and $C = C^H$, where $K^H$ denotes the conjugate transpose of $K$.

Sylvester equations play important roles in numerous applications such as matrix eigen-decompositions, control theory, model reduction, numerical solution of matrix differential Riccati equations, image processing, and many more, see for example [9, 1, 13, 14, 17, 25] and a large literature therein.

The Sylvester equation (1.1) is mathematically equivalent to the larger linear system of the form

$$\mathscr{A}\mathbf{x} = \mathbf{c}, \tag{1.2}$$

where $\mathscr{A} = I_m \otimes A + B^T \otimes I_n$ with $\otimes$ denoting the standard Kronecker product symbol, $\mathbf{x}$, $\mathbf{c}$ are two column-stacking vectors of the matrices $X$ and $C$, respectively. It is useful to treat the equation (1.2) as a general linear system in theoretical analysis, but impractical in numerical solution to the continuous Sylvester equation (1.1), because the equation (1.2) is costly to solve and can be ill-conditioned.

The common approaches to solving (1.1) are the Bartels-Stewart [12] and the Hessenberg-Schur [16, 17] methods, each of which needs to transform $A$ and $B$ into triangular or Hessenberg form by an orthogonal similarity transformation and then solving the resulting system directly by Gaussian elimination with partial pivoting. Those methods are usually referred to as the direct methods. The direct methods are mainly

---

*School of Mathematics and Statistics, Changsha University of Science and Technology, Changsha 410076, P. R. China (liuzhongyun@263.net)

†Centro de Matemática, Universidade do Minho, 4710-057 Braga, Portugal (Yulin Zhang: zhang@math.uminho.pt).

†A matrix $K$ is said to be positive definite or positive semi-definite if its Hermitian part, i.e., $\frac{1}{2}(K + K^H)$, is positive definite or positive semi-definite and a positive definite or positive semi-definite matrix is not necessarily Hermitian.

applicable for small and medium model problems but often too expensive to be practical for large sparse problems.

When $A$ and $B$ are large sparse matrices, an alternative for solving (1.1) is iterative methods. The Smith method [27] and the ADI iterative method [13, 15, 20, 22, 24, 26] are popular ones for solving large sparse Sylvester equations. Much more attention was paid to the case of the coefficient matrices $A$ and $B$ being either Hermitian positive definite matrices or $M$-matrices. However, little attention was focused on the case of the coefficient matrices $A$ and $B$ being non-Hermitian positive definite matrices. Recently, some authors applied the state-of-the-art iterative methods, such as the Hermitian and skew-Hermitian splitting (HSS) iteration [5], the positive-definite and skew-Hermitian splitting (PSS) iteration [8] and the normal and skew-Hermitian splitting (NSS) iteration [6] for solving non-Hermitian positive definite linear systems[¶], to develop HSS, PSS and NSS iteration solvers for (1.1). Of those methods, one also needs to solve two Sylvester equations with certain structural coefficient matrices at each inner iteration, those special structures allow the uses of the numerically stable Cholesky/Bunch-Parlett factorizations (as for the direct methods) or the short-term recurrence CG/GMRES (CGNE, CGNR) methods (as for the iteration methods). The resulting methods converge unconditionally and are efficient and robust numerically, see for example [1, 28, 29]. Nevertheless, the coefficient matrices may be dense (for instance, when the matrix $A$ is an upper Hessenberg matrix, $H$ and $S$ in the HSS splittings are still very dense), see [8]. This motivates us to further study the ADI-like iterations.

In this paper, we revisit the ADI iteration for solving (1.1). We first analyze its convergence and then derive an upper bound for its contraction factor. To reduce the computational complexity, we further establish an inexact variant of the ADI (IADI for short) iteration. The convergence of the IADI method is also analyzed in detail. Numerical experiments show that those methods are efficient and robust solvers and have a better performance than the HSS iteration solver for (1.1), and the IADI iteration is usually superior to the ADI iteration in terms of computation efficiency. Moreover, the IADI iteration as well as the ADI iteration outperform the IHSS iteration in terms of both the number of iterations and the computation efficiency.

The organization of this paper is as follows. After analyze the convergence rate of ADI iterative method for solving (1.1) in next section, we then establish the IADI iteration to improve the computing efficiency of the ADI iteration in section 3. Numerical experiments are illustrated in section 4 to show the effectiveness and robustness of our methods.

**2. The ADI Iteration.** Let us begin with some basic notations. For convenience, throughout this paper, we denote by $\mathbf{k} \in \mathbb{C}^{n^2}$ the column-stacking vector of the matrix $K \in \mathbb{C}^{n \times n}$, and we denote by $\lambda(K)$, $\rho(K)$, $\| K \|_2$ and $\| K \|_F$ the spectrum, spectral radius, the 2-norm, and the Frobenius norm of the matrix $K \in \mathbb{C}^{n \times n}$, respectively.

A matrix sequence $\{Y^{(k)}\}_{k=0}^{\infty} \subseteq \mathbb{C}^{n \times n}$ is said to be convergent to a matrix $Y \in \mathbb{C}^{n \times n}$ if the corresponding column-stacking vector sequence of $\{\mathbf{y}^{(k)}\}_{k=0}^{\infty} \subseteq \mathbb{C}^{n^2}$ of $\{Y^{(k)}\}_{k=0}^{\infty}$ is convergent to the corresponding column-stacking vector $\mathbf{y} \in \mathbb{C}^{n^2}$ of $Y$.

The classical ADI iterative method for solving (1.1) is as follows.

**The ADI iteration.** *Given an initial guess $X^{(0)}$, for $k = 0, 1, 2, \cdots$, until $\{X^{(k)}\}$ converges, compute*

$$\begin{cases} (\alpha I + A)X^{(k+\frac{1}{2})} = X^{(k)}(\alpha I - B) + C \\ X^{(k+1)}(\beta I + B) = (\beta I - A)X^{(k+\frac{1}{2})} + C, \end{cases} \tag{2.1}$$

*where $\alpha, \beta$ are given positive constants.*

---

[¶]In fact, such kind of methods are analogues to the classical ADI iteration introduced by Peaceman and Rachford for solving partial differential equations, also see [23].

Obviously, the equation (2.1) is equivalent to the following one

$$X^{(k+1)} = (\beta I - A)(\alpha I + A)^{-1} X^{(k)}(\alpha I - B)(\beta I + B)^{-1} + (\alpha + \beta)(\alpha I + A)^{-1} C(\beta I + B)^{-1}, \qquad (2.2)$$

which becomes the Smith iteration [27] when $\alpha = \beta$. Therefore, we can think of the ADI iteration as an accelerated or a generalized version of the Smith iteration.

We remark here that the two half-steps at each iteration (2.2) for solving Sylvester equation (1.1) need to solve two linear subsystems. Because the coefficient matrices are positive (semi-) definite, we can suitably choose the shifts in such a way that those matrices have reasonably good diagonally dominant property such that their incomplete factorizations are existent, stable, and accurate, as mentioned in [10]; or good preconditioners make the GMRES converge fast, for some effective preconditioner for Hermitian and non-Hermitian positive definite matrices, see [2, 3, 4].

In matrix-vector form, the ADI iteration (2.1) can be equivalently rewritten as

$$\mathbf{x}^{(k+1)} = \mathscr{H}(\alpha, \beta)\mathbf{x}^{(k)} + \mathscr{G}\mathbf{c}, \qquad (2.3)$$

where

$$\begin{cases} \mathscr{H}(\alpha, \beta) = [(\beta I + B)^{-T}(\alpha I - B)^T] \otimes [(\beta I - A)(\alpha I + A)^{-1}], \\ \mathscr{G} = (\alpha + \beta)[(\beta I + B)^{-T} \otimes (\alpha I + A)^{-1}], \end{cases}$$

and $\mathbf{x} \in \mathbb{C}^{n^2}$, $\mathbf{c} \in \mathbb{C}^{n^2}$, and $\mathscr{H}(\alpha, \beta)$ is the iteration matrix of (2.3).

Before giving the convergence theorem of the ADI iteration, we recall the following known results.

LEMMA 2.1. [19] *For any* $A, B \in \mathbb{C}^{n \times n}$, $\rho(AB) = \rho(BA)$ *and* $\rho(A \otimes B) = \rho(A) \cdot \rho(B)$.

LEMMA 2.2. *Let* $K \in \mathbb{C}^{n \times n}$, $\lambda_j \in \lambda(K)$. *Then* $\rho[(\alpha I + K)^{-1}(\beta I - K)] = \max\limits_{1 \le j \le n} \left| \frac{\beta - \lambda_j}{\alpha + \lambda_j} \right|$.

Now, we can give the convergence theorem of the ADI iteration (2.1).

THEOREM 2.3. *Let* $A \in \mathbb{C}^{n \times n}$ *be positive definite and* $B \in \mathbb{C}^{n \times n}$ *be positive semi-definite,* $\alpha, \beta$ *be two positive constants. For* $j = 1, \cdots, n$, *let* $\lambda_j = \lambda_j' + i\lambda_j''$ *and* $\mu_j = \mu_j' + i\mu_j''$ *be the eigenvalues of the matrices* $A$ *and* $B$, *respectively, where* $\lambda_j'$, $\lambda_j''$ *and* $\mu_j'$, $\mu_j''$ *are the real and pure imaginary parts of the eigenvalues* $\lambda_j$ *and* $\mu_j$. *Let* $\tau = \frac{\alpha + \beta}{2}$ *and* $\Delta = \frac{\alpha - \beta}{2}$. *Then the iterative sequence* $\{X^{(k)}\}$ *determined by (2.1) converges to the exact solution* $X^*$ *of (1.1), provided that* $-\lambda_{\min}' < \Delta < \mu_{\min}'$, *where* $\lambda_{\min}'$ *and* $\mu_{\min}'$ *denote the lower bounds of the real parts of the eigenvalues of the matrices* $A$ *and* $B$, *respectively.*

*Proof.* From (2.3) and by Lemma 2.1 and Lemma 2.2, we have

$$\rho[\mathscr{H}(\alpha, \beta)] = \rho\{[(\beta I + B)^{-T}(\alpha I - B)^T] \otimes [(\beta I - A)(\alpha I + A)^{-1}]\}$$

$$= \max_{\lambda_j \in \lambda(A), \ \mu_j \in \lambda(B)} \left| \frac{\alpha - \mu_j}{\beta + \mu_j} \cdot \frac{\beta - \lambda_j}{\alpha + \lambda_j} \right|.$$

Denoting

$$\phi_1(\alpha, \beta) = \max_{\lambda_j \in \lambda(A)} \left| \frac{\beta - \lambda_j}{\alpha + \lambda_j} \right| \quad \text{and} \quad \phi_2(\alpha, \beta) = \max_{\mu_j \in \lambda(B)} \left| \frac{\alpha - \mu_j}{\beta + \mu_j} \right|, \qquad (2.4)$$

we get

$$\rho[\mathscr{H}(\alpha, \beta)] \le \phi_1(\alpha, \beta) \cdot \phi_2(\alpha, \beta)$$

$$= \max_{\lambda_j = \lambda_j' + i\lambda_j'' \in \lambda(A)} \sqrt{\frac{(\beta - \lambda_j')^2 + (\lambda_j'')^2}{(\alpha + \lambda_j')^2 + (\lambda_j'')^2}} \cdot \max_{\mu_j = \mu_j' + i\mu_j'' \in \lambda(B)} \sqrt{\frac{(\alpha - \mu_j')^2 + (\mu_j'')^2}{(\beta + \mu_j')^2 + (\mu_j'')^2}}.$$

Setting $\tilde{\lambda}_j = (\lambda'_j + \Delta) + i\lambda''_j = \tilde{\lambda}'_j + i\lambda''_j$, $\tilde{\mu}_j = (\mu'_j - \Delta) + i\mu''_j = \tilde{\mu}'_j + i\mu''_j$, the functions $\phi_1(\alpha, \beta)$ and $\phi_2(\alpha, \beta)$ in (2.4) become

$$\phi_1(\tau, \Delta) = \max_{\tilde{\lambda}_j = \tilde{\lambda}'_j + i\lambda''_j} \sqrt{\frac{(\tau - \tilde{\lambda}'_j)^2 + (\lambda''_j)^2}{(\tau + \tilde{\lambda}'_j)^2 + (\lambda''_j)^2}} \quad \text{and} \quad \phi_2(\tau, \Delta) = \max_{\tilde{\mu}_j = \tilde{\mu}'_j + i\mu''_j} \sqrt{\frac{(\tau - \tilde{\mu}'_j)^2 + (\mu''_j)^2}{(\tau + \tilde{\mu}'_j)^2 + (\mu''_j)^2}} \tag{2.5}$$

we then have

$$\rho[\mathscr{H}(\alpha, \beta)] \leq \phi_1(\tau, \Delta) \cdot \phi_2(\tau, \Delta) \equiv \phi(\tau, \Delta).$$

Due to $-\lambda'_{\min} < \Delta < \mu'_{\min}$, we have that $\tilde{\lambda}'_j > 0$ and $\tilde{\mu}'_j > 0$, which imply that $\phi_1(\tau, \Delta) < 1$, $\phi_2(\tau, \Delta) < 1$ and so $\rho[\mathscr{H}(\alpha, \beta)] \leq \phi(\tau, \Delta) < 1$. Thus the proof is complete. □

In the following, we will suggest a method to determine the $\alpha$ and $\beta$ in (2.1). Consider $\Delta$ as a parameter in $\phi(\tau, \Delta), \phi_1(\tau, \Delta)$ and $\phi_2(\tau, \Delta)$. Note that for any fixed $\Delta \in (-\lambda'_{\min}, \mu'_{\min})$,

$$\min_\tau \phi(\tau, \Delta) \geq \min_\tau \phi_1(\tau, \Delta) \cdot \min_\tau \phi_2(\tau, \Delta)$$

and the equality holds if and only if for a certain $\Delta^*$ there is a corresponding $\tau^*$ that minimizes $\phi_1(\tau, \Delta^*)$ and $\phi_2(\tau, \Delta^*)$ (as well as $\phi(\tau, \Delta^*)$) simultaneously. If $(\tau^*, \Delta^*)$ can be determined in this way and in addition it satisfies $\tau^* > \Delta^*$, we set $\alpha^* = \tau^* + \Delta^*$ and $\beta^* = \tau^* - \Delta^*$.

Using an argument similar to the Theorem 2.2 in [6, 7], the minimizers $\tau_1^*, \tau_2^*$ for solving $\min_\tau \phi_1(\tau, \Delta)$ and $\min_\tau \phi_2(\tau, \Delta)$ for a fixed $\Delta$, respectively, are given by

$$\tau_1^* = \begin{cases} \sqrt{(\lambda'_{\min} + \Delta)(\lambda'_{\max} + \Delta) - {\lambda''_{\max}}^2}, & \text{for } \lambda''_{\max} < \theta_1 \\ \sqrt{(\lambda'_{\min} + \Delta)^2 + {\lambda''_{\max}}^2}, & \text{for } \lambda''_{\max} \geq \theta_1, \end{cases} \tag{2.6}$$

and

$$\tau_2^* = \begin{cases} \sqrt{(\mu'_{\min} - \Delta)(\mu'_{\max} - \Delta) - {\mu''_{\max}}^2}, & \text{for } \mu''_{\max} < \theta_2, \\ \sqrt{(\mu'_{\min} - \Delta)^2 + {\mu''_{\max}}^2}, & \text{for } \mu''_{\max} \geq \theta_2. \end{cases} \tag{2.7}$$

where

$$\theta_1 = \sqrt{\frac{(\lambda'_{\min} + \Delta)(\lambda'_{\max} - \lambda'_{\min})}{2}} \quad \text{and} \quad \theta_2 = \sqrt{\frac{(\mu'_{\min} - \Delta)(\mu'_{\max} - \mu'_{\min})}{2}}.$$

By setting $\tau_1^* = \tau_2^*$ and solving it for $\Delta$, we have

$$\begin{cases} \Delta_1 = \frac{\mu'_{\min}\mu'_{\max} - {\mu''_{\max}}^2 + {\lambda''_{\max}}^2 - \lambda'_{\min}\lambda'_{\max}}{\lambda'_{\min} + \lambda'_{\max} + \mu'_{\min} + \mu'_{\max}}, & \text{when } \lambda''_{\max} < \theta_1 \text{ and } \mu''_{\max} < \theta_2 \\[2mm] \Delta_2 = \frac{\mu'_{\min}\mu'_{\max} - {\mu''_{\max}}^2 - {\lambda''_{\max}}^2 - {\lambda'_{\min}}^2}{2\lambda'_{\min} + \mu'_{\min} + \mu'_{\max}}, & \text{when } \lambda''_{\max} < \theta_1 \text{ and } \mu''_{\max} \geq \theta_2 \\[2mm] \Delta_3 = \frac{{\mu'_{\min}}^2 + {\mu''_{\max}}^2 - {\lambda'_{\min}}^2 - {\lambda''_{\max}}^2}{2\lambda'_{\min} + 2\mu'_{\min}}, & \text{when } \lambda''_{\max} \geq \theta_1 \text{ and } \mu''_{\max} < \theta_2 \\[2mm] \Delta_4 = \frac{{\mu'_{\min}}^2 + {\mu''_{\max}}^2 - \lambda'_{\min}\lambda'_{\max} + {\lambda''_{\max}}^2}{\lambda'_{\min} + \lambda'_{\max} + 2\mu'_{\min}}, & \text{when } \lambda''_{\max} \geq \theta_1 \text{ and } \mu''_{\max} \geq \theta_2. \end{cases} \tag{2.8}$$

For each $\Delta_i$, we calculate the corresponding $\theta_1, \theta_2$ values denoted by $\theta_1^{(i)}, \theta_2^{(i)}, i = 1, ..., 4$, respectively. We then determine $\Delta^*$ with the following rule

$$\Delta^* = \begin{cases} \Delta_1, & \text{when} - \lambda'_{\min} < \Delta_1 < \mu'_{\min} \text{ and } \lambda''_{\max} < \theta_1^{(1)} \text{ and } \mu''_{\max} < \theta_2^{(1)} \\[2ex] \Delta_2, & \text{when} - \lambda'_{\min} < \Delta_2 < \mu'_{\min} \text{ and } \lambda''_{\max} < \theta_1^{(2)} \text{ and } \mu''_{\max} \geq \theta_2^{(2)} \\[2ex] \Delta_3, & \text{when} - \lambda'_{\min} < \Delta_3 < \mu'_{\min} \text{ and } \lambda''_{\max} \geq \theta_1^{(3)} \text{ and } \mu''_{\max} < \theta_2^{(3)} \\[2ex] \Delta_4, & \text{when} - \lambda'_{\min} < \Delta_4 < \mu'_{\min} \text{ and } \lambda''_{\max} \geq \theta_1^{(4)} \text{ and } \mu''_{\max} \geq \theta_2^{(4)}. \end{cases} \quad (2.9)$$

Once such a $\Delta^*$ is determined, we determine $\tau^* = \tau_1^* = \tau_2^*$ by using the formulas either in (2.6) or (2.7). If there are several solutions $(\tau^*, \Delta^*)$, then we pick up a pair with the smallest $\phi(\tau^*, \Delta^*)$.

If no $\Delta^*$ can be determined, or it happens that $\tau^* \leq \Delta^*$ or $\Delta^* \leq -\tau^*$ for all determined pairs, the above procedure fails to determine $(\alpha^*, \beta^*)$. When this happens, we consider $\Delta = 0$ (i.e., $\alpha = \beta$).

In this case, ADI iteration (2.1) reduces to Smith iteration, $\rho[\mathscr{H}(\alpha)]$ is bounded by

$$\hat{\phi}(\alpha) = \max_{\gamma' + i\gamma'' \in \Omega} \frac{(\alpha - \gamma')^2 + (\gamma'')^2}{(\alpha + \gamma')^2 + (\gamma'')^2},$$

where $\Omega = [\gamma'_{min}, \gamma'_{max}] \times i[\gamma''_{min}, \gamma''_{max}]$ with $\gamma'_{min}$ and $\gamma'_{max}$ denoting the lower and the upper bounds of the real part of the eigenvalues of the matrices $A$ and $B$, and $\gamma''_{min}$ and $\gamma''_{max}$ denoting the lower and the upper bounds of the absolute values of the imaginary part of the eigenvalues of the matrices $A$ and $B$. The parameter $\alpha^*$ is chosen such that the above estimate can be minimized. This fact is precisely stated as the following theorem.

LEMMA 2.4. *The minimizer of $\hat{\phi}(\alpha)$ over all positive $\alpha$ is attained at*

$$\alpha^* = \arg\min_\alpha \hat{\phi}(\alpha) = \begin{cases} \sqrt{\gamma'_{\min}\gamma'_{\max} - \gamma''^2_{\max}}, & \text{for } \gamma''_{\max} < \sqrt{\frac{\gamma'_{\min}(\gamma'_{\max} - \gamma'_{\min})}{2}} \\[2ex] \sqrt{\gamma'^2_{\min} + \gamma''^2_{\max}}, & \text{for } \gamma''_{\max} \geq \sqrt{\frac{\gamma'_{\min}(\gamma'_{\max} - \gamma'_{\min})}{2}}, \end{cases} \quad (2.10)$$

*and the corresponding minimum value is equal to*

$$\hat{\phi}(\alpha^*) = \begin{cases} \dfrac{\gamma'_{\min} + \gamma'_{\max} - 2\sqrt{\gamma'_{\min}\gamma'_{\max} - \gamma''^2_{\max}}}{\gamma'_{\min} + \gamma'_{\max} + 2\sqrt{\gamma'_{\min}\gamma'_{\max} - \gamma''^2_{\max}}}, & \text{for } \gamma''_{\max} < \sqrt{\frac{\gamma'_{\min}(\gamma'_{\max} - \gamma'_{\min})}{2}} \\[3ex] \dfrac{\sqrt{\gamma'^2_{\min} + \gamma''^2_{\max}} - \gamma'_{\min}}{\sqrt{\gamma'^2_{\min} + \gamma''^2_{\max}} + \gamma'_{\min}}, & \text{for } \gamma''_{\max} \geq \sqrt{\frac{\gamma'_{\min}(\gamma'_{\max} - \gamma'_{\min})}{2}}. \end{cases}$$

The proof is a verbatim of one of Corollary 2.3 in [5] and therefore omitted.

REMARK 2.5. *We remark that the condition $-\lambda'_{\min} < \Delta < \mu'_{\min}$ is sufficient but not necessary. That is to say that even if $-\lambda'_{\min} < \Delta < \mu'_{\min}$ does not hold, we cannot guarantee that the upper bound $\phi$ is less than 1, but it does not mean the spectral radius $\rho[\mathscr{H}(\alpha, \beta)]$ cannot be less than 1.*

In fact, if the imaginary parts of the eigenvalues are zero, for example, when $A$ and $B$ are both Hermitian positive definite matrices, then condition $-\lambda'_{\min} < \Delta < \mu'_{\min}$ always holds. When $A$ and $B$ are both non-Hermitian positive definite matrices, we cannot guarantee that $-\lambda'_{\min} < \Delta < \mu'_{\min}$. From (2.8), however, we observe that if one of the following cases:

(1) $\lambda''_{\max}$ and $\mu''_{\max}$ are small enough,
(2) $||\lambda''_{\max}| - |\mu''_{\max}||$ is sufficient small,
(3) $\lambda'_{\max}$ (and $\mu'_{\max}$) is much larger than $\lambda''_{\max}$ (and $\mu''_{\max}$)

holds, then we can find a $\Delta$ such that $-\lambda'_{\min} < \Delta < \mu'_{\min}$ holds. This phenomenon appears in our numerical tests. Therefore, we have the following corollary.

COROLLARY 2.6.  *Under the assumption of Theorem 2.3, the iterative sequence $\{X^{(k)}\}$ determined by (2.1) converges to the exact solution $X^*$ of (1.1), if $\lambda''_{\max}$ and $\mu''_{\max}$ are small enough, $||\lambda''_{\max}| - |\mu''_{\max}||$ is sufficient small, or $\lambda'_{\max}$ (and $\mu'_{\max}$) is much larger than $\lambda''_{\max}$ (and $\mu''_{\max}$).*

COROLLARY 2.7.  *If $\alpha = \beta$, i.e., $\Delta = 0$, then the iteration (2.1) converges to the exact solution unconditionally for all $\alpha > 0$.*

If $B = A^T$ in (1.1), then $\lambda(B) = \lambda(A)$. From (2.8), we therefore have $\Delta = 0$. In this case we obtain the following result.

COROLLARY 2.8.  *If $B = A^T$ in (1.1), then the iteration (2.1) converges to the exact solution unconditionally for all $\alpha > 0$.*

**3. The IADI Iteration.** The two half-steps at each step of the ADI iteration for solving continuous Sylvester equation (1.1) need to solve two matrix equations like

$$(\alpha I + A)X = Y_1 \quad \text{and} \quad X(\beta I + B) = Y_2, \tag{3.1}$$

where $Y_1$ and $Y_2$ are known.

This may be very costly and impractical in actual implementations. To further improve the computational efficiency of the ADI iteration, we can solve the two subproblems in (3.1) inexactly by employing some state of the art iterative methods such as GMRES, which results in the basic framework of the IADI iterative method for solving (1.1).

Given an initial guess $X^{(0)} \in \mathbb{C}^{n \times n}$, for $k = 1, 2, \cdots$, until $\{X^{(k)}\}_{k=0}^{\infty} \subseteq \mathbb{C}^{n \times n}$ satisfies the stopping criterion, solve $X^{(k+\frac{1}{2})}$ approximately from

$$(\alpha I + A)X^{(k+\frac{1}{2})} \approx X^{(k)}(\alpha I - B) + C,$$

by employing an inner iteration (e.g., the GMRES) with $X^{(k)}$ as the initial guess; then solve $X^{(k+1)}$ approximately from

$$X^{(k+1)}(\beta I + B) \approx (\beta I - A)X^{(k+\frac{1}{2})} + C,$$

by employing an inner iteration (e.g., GMRES) with $X^{(k+\frac{1}{2})}$ as the initial guess, where $\alpha, \beta$ are given positive constants.

To simplify numerical implementation and convergence analysis, based on the residual-updating form [11] we may rewrite the above IADI iteration as the following equivalent scheme.

**The IADI Iteration.**  *Given an initial guess $X^{(0)} \in \mathbb{C}^{n \times n}$, for $k = 1, 2, \cdots$, until $\{X^{(k)}\}_{k=0}^{\infty} \subseteq \mathbb{C}^{n \times n}$ converges:*

  *1. approximate the solution of*

$$(\alpha I + A)Z^{(k)} = R^{(k)},$$

*with $R^{(k)} = C - AX^{(k)} - X^{(k)}B$, by iterating until $Z^{(k)}$ is such that the residual*

$$P^{(k)} = R^{(k)} - (\alpha I + A)Z^{(k)}$$

satisfies

$$\| P^{(k)} \|_F \leq \varepsilon_k \| R^{(k)} \|_F,$$

and then compute

$$X^{(k+\frac{1}{2})} = X^{(k)} + Z^{(k)};$$

2. approximate the solution of

$$Z^{(k+\frac{1}{2})}(\beta I + B) = R^{(k+\frac{1}{2})},$$

with $R^{(k+\frac{1}{2})} = C - AX^{(k+\frac{1}{2})} - X^{(k+\frac{1}{2})}B$, by iterating until $Z^{(k+\frac{1}{2})}$ is such that the residual

$$Q^{(k+\frac{1}{2})} = R^{(k+\frac{1}{2})} - Z^{(k+\frac{1}{2})}(\beta I + B)$$

satisfies

$$\| Q^{(k+\frac{1}{2})} \|_F \leq \eta_k \| R^{(k+\frac{1}{2})} \|_F,$$

and then compute

$$X^{(k+1)} = X^{(k+\frac{1}{2})} + Z^{(k+\frac{1}{2})},$$

where $\{\varepsilon_k\}$ and $\{\eta_k\}$ are prescribed tolerances used to control the accuracies of the inner iterations.

In order to analyze the convergence of the above IADI iteration, we need to recall the following convergence theorem of iterative solution to a general linear system $Ax = b$ by an inexact two-step splitting iterative method. Let $||v||$ (for all $v \in \mathbb{C}^n$) be a general vector norm and $M$ be a nonsingular matrix. We define a new vector norm by $|||v|||_M = ||Mv||$ (for all $v \in \mathbb{C}^n$), then its induced matrix norm is $|||K|||_M = ||MKM^{-1}||$ (for all $K \in \mathbb{C}^{n \times n}$ ).

LEMMA 3.1. [5, Theorem 3.1] Let $A \in \mathbb{C}^{n \times n}$ and $A = M_i - N_i$ $(i = 1, 2)$ be two splitings of the matrix $A$. If $\{\tilde{x}^{(k)}\}$ is an iterative sequence defined as

$$\tilde{x}^{(k+\frac{1}{2})} = \tilde{x}^{(k)} + \tilde{z}^{(k)}, \quad with \quad M_1\tilde{z}^{(k)} = \tilde{r}^{(k)} + \tilde{p}^{(k)},$$

satisfying $\frac{||\tilde{p}^{(k)}||}{||\tilde{r}^{(k)}||} \leq \varepsilon_k$, where $\tilde{r}^{(k)} = b - A\tilde{x}^{(k)}$, and

$$\tilde{x}^{(k+1)} = \tilde{x}^{(k+\frac{1}{2})} + \tilde{z}^{(k+\frac{1}{2})}, \quad with \quad M_2\tilde{z}^{(k+\frac{1}{2})} = \tilde{r}^{(k+\frac{1}{2})} + \tilde{q}^{(k+\frac{1}{2})},$$

satisfying $\frac{||\tilde{q}^{(k+\frac{1}{2})}||}{||\tilde{r}^{(k+\frac{1}{2})}||} \leq \eta_k$, where $\tilde{r}^{(k+\frac{1}{2})} = b - A\tilde{x}^{(k+\frac{1}{2})}$, then $\{\tilde{x}^{(k)}\}$ is of the form

$$\tilde{x}^{(k+1)} = M_2^{-1}N_2M_1^{-1}N_1\tilde{x}^{(k)} + M_2^{-1}(I + N_2M_1^{-1})b + M_2^{-1}(N_2M_1^{-1}\tilde{p}^{(k)} + \tilde{q}^{(k+\frac{1}{2})}).$$

Moreover, if $x^* \in \mathbb{C}^n$ is the exact solution of the linear system $Ax = b$, then we have

$$|||\tilde{x}^{(k+1)} - x^*|||_{M_2} \leq \tau_k|||\tilde{x}^{(k)} - x^*|||_{M_2}, \quad k = 0, 1, 2, \cdots,$$

where $\tau_k = \sigma + \mu\theta\varepsilon_k + \theta(\rho + \theta\nu\varepsilon_k)\eta_k$ with

$$\sigma = ||N_2M_1^{-1}N_1M_2^{-1}||, \qquad \rho = ||M_2M_1^{-1}N_1M_2^{-1}||, \qquad \mu = ||N_2M_1^{-1}||,$$
$$\theta = ||AM_2^{-1}||, \qquad \nu = ||M_2M_1^{-1}||.$$

*In particular, if*

$$\tau_{max} \equiv \sigma + \mu\theta\varepsilon_{max} + \theta(\rho + \theta\nu\varepsilon_{max})\eta_{max} < 1,$$

*then the iterative sequence $\{\tilde{x}^{(k)}\}$ converges to $x^* \in \mathbb{C}^n$, where $\varepsilon_{\max} = \max\limits_{k} \varepsilon_k$ and $\eta_{\max} = \max\limits_{k} \eta_k$.*

Now we can demonstrate the following convergence result concerning the above IADI iterative method.

THEOREM 3.2. *Let $\{X^{(k)}\}_{k=0}^{\infty} \subseteq \mathbb{C}^{n \times n}$ be an iteration sequence generated by the IADI iterative method and let $X^* \in \mathbb{C}^{n \times n}$ be the exact solution of* (1.1). *Under the assumption of Theorem 2.3, we have that*

$$\| X^{(k+1)} - X^* \|_B \,^{\dagger\dagger} \leq \tau_k \| X^{(k)} - X^* \|_B \tag{3.2}$$

*where*

$$\tau_k = \sigma + \mu\theta\varepsilon_k + \theta(\sigma + \theta\nu\varepsilon_k)\eta_k \tag{3.3}$$

*with*

$$\sigma = \| [(\alpha I - B)^T (\beta I + B)^{-T}] \otimes [(\beta I - A)(\alpha I + A)^{-1}] \|_2,$$
$$\mu = \| I \otimes (\beta I - A)(\alpha I + A)^{-1} \|_2,$$
$$\theta = \| (\beta I + B)^{-T} \otimes A + B^T (\beta I + B)^{-T} \otimes I \|_2,$$
$$\nu = \| (\beta I + B)^T \otimes (\alpha I + A)^{-1} \|_2 .$$

*In particular, if*

$$\tau_{\max} \equiv \sigma + \mu\theta\varepsilon_{\max} + \theta(\sigma + \theta\nu\varepsilon_{\max})\eta_{\max} < 1, \tag{3.4}$$

*then the iteration sequence $\{X^{(k)}\}_{k=0}^{\infty}$ converges to $X^*$, where $\varepsilon_{\max} = \max\limits_{k} \varepsilon_k$ and $\eta_{\max} = \max\limits_{k} \eta_k$.*

*Proof.* Denoting $\mathscr{M}_1 = I \otimes (\alpha I + A), \mathscr{M}_2 = (\beta I + B)^T \otimes I, \mathscr{N}_1 = (\alpha I - B)^T \otimes I$ and $\mathscr{N}_2 = I \otimes (\beta I - A)$, we have that $\mathscr{A} = \mathscr{M}_i - \mathscr{N}_i$ $(i = 1, 2)$ are two splittings of $\mathscr{A}$. Again, by making use of the kronecker product, we can rewrite the above-described IADI iteration in the following matrix-vector form:

$$\mathscr{M}_1 \mathbf{z}^{(k)} = \mathbf{r}^{(k)}, \quad \mathbf{x}^{(k+\frac{1}{2})} = \mathbf{x}^{(k)} + \mathbf{z}^{(k)} \tag{3.5}$$

with $\mathbf{r}^{(k)} = \mathbf{c} - \mathscr{A}\mathbf{x}^{(k)}$, where $\mathbf{z}^{(k)}$ is the approximate solution such that the residual

$$\mathbf{p}^{(k)} = \mathbf{r}^{(k)} - \mathscr{M}_1 \mathbf{z}^{(k)}$$

satisfies $\| \mathbf{p}^{(k)} \|_2 \leq \varepsilon_k \| \mathbf{r}^{(k)} \|_2$, and

$$\mathscr{M}_2 \mathbf{z}^{(k+\frac{1}{2})} = \mathbf{r}^{(k+\frac{1}{2})}, \quad \mathbf{x}^{(k+1)} = \mathbf{x}^{(k+\frac{1}{2})} + \mathbf{z}^{(k+\frac{1}{2})} \tag{3.6}$$

with $\mathbf{r}^{(k+\frac{1}{2})} = \mathbf{c} - \mathscr{A}\mathbf{x}^{(k+\frac{1}{2})}$, where $\mathbf{z}^{(k+\frac{1}{2})}$ is the approximate solution such that the residual

$$\mathbf{q}^{(k+\frac{1}{2})} = \mathbf{r}^{(k+\frac{1}{2})} - \mathscr{M}_2 \mathbf{z}^{(k+\frac{1}{2})}$$

satisfies

$$\| \mathbf{q}^{(k+\frac{1}{2})} \|_2 \leq \eta_k \| \mathbf{r}^{(k+\frac{1}{2})} \|_2 .$$

By Lemma 3.1, we can easily obtain

$$|||\mathbf{x}^{(k+1)} - \mathbf{x}^*|||_{\mathscr{M}_2} \leq \tau_k |||\mathbf{x}^{(k)} - \mathbf{x}^*|||_{\mathscr{M}_2}. \tag{3.7}$$

---

$^{\dagger\dagger}$For any $K \in \mathbb{C}^{n \times n}$, we define its matrix norm by $\| K \|_B = \| K(\beta I + B) \|_F$.

Now, taking $|||\mathbf{y}|||_{\mathscr{M}_2} =|| \mathscr{M}_2\mathbf{y} ||_2$, then we have

$$|||\mathbf{y}|||_{\mathscr{M}_2} =|| \mathscr{M}_2\mathbf{y} ||_2=|| [(\beta I + B)^T \otimes I]\mathbf{y} ||_2=|| Y(\beta I + B) ||_F=|| Y ||_B .$$

Hence, we can equivalently rewrite the inequality (3.7) as

$$|| X^{(k+1)} - X^* ||_B \leq \tau_k || X^{(k)} - X^* ||_B .$$

Thus we complete the proof of this theorem. □

According to Theorem 3.2, we want to choose tolerances $\{\varepsilon_k\}$ and $\{\eta_k\}$ so that the computational work of the IADI iterative method is minimized. In fact, the tolerances $\{\varepsilon_k\}$ and $\{\eta_k\}$ are not required to approach zero as $k$ increases in order to ensure the convergence of the IADI iteration but are required to approach zero in order to asymptotically recover the original convergence rate (cf. Theorem 2.3) of the ADI iteration. How to arrive at a tradeoff between the computational complexity and the convergence rate is a difficult optimal problem, it deserves further in-depth study.

The following theorem presents one possible way of choosing the tolerances $\{\varepsilon_k\}$ and $\{\eta_k\}$ such that the original convergence rate of the ADI iteration can be asymptotically recovered.

THEOREM 3.3. *Let the assumptions in Theorem 3.2 be satisfied. Suppose that both $\{\psi_1(k)\}$ and $\{\psi_2(k)\}$ are nondecreasing and positive sequences satisfying $\psi_1(k) \geq 1$, $\psi_2(k) \geq 1$, and $\lim\limits_{k\to\infty} \sup \psi_1(k) = \lim\limits_{k\to\infty} \sup \psi_2(k) = +\infty$, and that both $\delta_1$ and $\delta_2$ are real constants on the interval $(0,1)$ satisfying*

$$\varepsilon_k \leq t_1\delta_1^{\psi_1(k)} \quad and \quad \eta_k \leq t_2\delta_2^{\psi_2(k)}, \tag{3.8}$$

*where $t_1$ and $t_2$ are nonnegative constants. Then we have*

$$|| X^{(k+1)} - X^* ||_B \leq (\sqrt{\sigma} + \varphi\theta\delta^{\psi(k)})^2 || X^{(k)} - X^* ||_B,$$

*where*

$$\psi(k) = \min\{\psi_1(k),\ \psi_2(k)\}, \quad \delta = \max\{\delta_1,\ \delta_2\},$$

*and*

$$\varphi = \max\{\sqrt{t_1 t_2 \nu},\ \frac{1}{2\sqrt{\sigma}}(t_1\mu + t_2\sigma)\},$$

*In particular, we have*

$$\lim\limits_{k\to\infty} \sup \frac{|| X^{(k+1)} - X^* ||_B}{|| X^{(k)} - X^* ||_B} \leq \sigma.$$

*i.e. the convergence rate of the IADI iterative method is asymptotically the same as that of the ADI iterative method.*

*Proof.* The proof is a verbatim of Theorems 3.3 and 3.4 in [5] and thus omitted. □

**4. Numerical examples.** In this section, we use some examples to illustrate the effectiveness of ADI and IADI iterations for solving the Sylvester equation(1.1).

In actual computations, all iterations are started from the zero matrix, performed in MATLAB with machine precision $10^{-16}$, and stopped if the norm of the current residual matrix satisfies $|| R^{(k)} ||_F\ /\ || R^{(0)} ||_F \leq 10^{-6}$, where, following the definition, $R^{(k)} = C - AX^{(k)} - X^{(k)}B$ is the residual matrix of the $k$-th ADI iteration.

For convenience, throughout our numerical experiments, we denote by **IT** the number of the iteration steps, by **CPU** the computing time (in seconds), and by $\alpha_{\exp}$ and $\beta_{\exp}$ the experimentally found local optimal values[‡] of the iteration parameters of the ADI iteration, respectively.

EXAMPLE 4.1. *We consider the Sylvester equation* (1.1) *and the matrices*

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

*where* $M, N \in \mathbb{R}^{n \times n}$ *are the tridiagonal matrices given by*

$$M = tridiag(-1, 2, -1) \quad and \quad N = tridiag(0.5, 0, -0.5).$$

Such a class of problems arise frequently in the preconditioned Krylov subspace iterative methods, see [1] and reference therein. For comparison, we also test the HSS solver. The numerical results obtained by the ADI method and the HSS solver are listed in Table 4.1 with $r = 1$, Table 4.2 with $r = 0.1$, Table 4.3 with $r = 0.01$, where let $\beta_{\exp} = \alpha_{\exp}$, and $\omega_{\exp}$ is the numerically found optimal value of the shift parameter of the HSS iterations in [1]. We can see that both the number of iteration steps and the runtime by ADI iteration are much less than those by HSS iteration in all cases.

In Fig.4.1, we depict the convergence behavior of the ADI and HSS iterations for Example 4.1, which are denoted by $+++$ and $\circ\circ\circ$ curves, respectively. From Fig.4.1, we can see that the ADI iteration has a better convergence behavior than the HSS iteration. We can also observe that the HSS is not sensitive to the nonsymmetric part, as mentioned in [1], where the author pointed out that the convergence depends only on the spectrums of the Hermitian parts (symmetric parts in our example).

Table 4.1 IT and CPU for ADI and HSS with r=1

|     | ADI | | | HSS | | |
| --- | --- | --- | --- | --- | --- | --- |
| $n$ | $\alpha_{\exp}$ | IT | CPU | $\omega_{\exp}$ | IT | CPU |
| 32 | 1.20 | 12 | 0.001 | 0.95 | 27 | 0.234 |
| 64 | 0.88 | 17 | 0.004 | 0.81 | 44 | 1.614 |
| 128 | 0.62 | 24 | 0.015 | 0.62 | 93 | 9.841 |
| 256 | 0.51 | 32 | 0.250 | 0.51 | 203 | 60.708 |

Table 4.2 IT and CPU for ADI and HSS with r=0.1

|     | ADI | | | HSS | | |
| --- | --- | --- | --- | --- | --- | --- |
| $n$ | $\alpha_{\exp}$ | IT | CPU | $\omega_{\exp}$ | IT | CPU |
| 32 | 0.74 | 18 | 0.002 | 0.4 | 48 | 0.488 |
| 64 | 0.43 | 31 | 0.010 | 0.23 | 92 | 2.943 |
| 128 | 0.27 | 50 | 0.063 | 0.13 | 177 | 20.721 |
| 256 | 0.18 | 74 | 0.687 | 0.09 | 274 | 161.132 |

Table 4.3 IT and CPU for ADI and HSS with r=0.01

|     | ADI | | | HSS | | |
| --- | --- | --- | --- | --- | --- | --- |
| $n$ | $\alpha_{\exp}$ | IT | CPU | $\omega_{\exp}$ | IT | CPU |
| 32 | 0.75 | 18 | 0.005 | 0.4 | 27 | 0.390 |
| 64 | 0.42 | 32 | 0.010 | 0.17 | 44 | 4.299 |
| 128 | 0.25 | 55 | 0.062 | 0.09 | 93 | 36.051 |
| 256 | 0.15 | 92 | 0.827 | 0.05 | 203 | 429.960 |

---

[‡]The scalars $\alpha_{\exp}$ and $\beta_{\exp}$ are obtained by searching the optimal values of the iteration parameters for the ADI iteration in two intervals $(\alpha - 1, \alpha + 1)$ and $(\beta - 1, \beta + 1)$ with stepsize 0.1, respectively, where $\alpha$ and $\beta$ are the approximate values obtained from $\Delta$ and $\tau$ defined as in (2.8) and (2.6), in which $\lambda_{min}$ and $\mu_{min}$ are approximately computed by employing the inverse iteration, and $\lambda_{max}$ and $\mu_{max}$ are roughly estimated by using the power iteration.
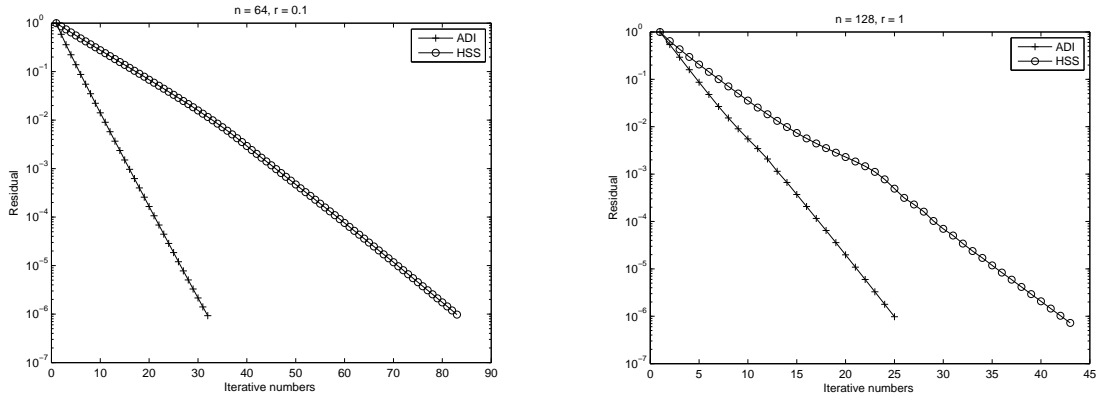
Fig. 4.1 Convergence curves for ADI and HSS for $n = 64, r = 0.1$ (left) and $n = 128, r = 1$ (right) in Example 4.1.

EXAMPLE 4.2. *We consider the Sylvester equation* (1.1) *and the matrices*

$$\begin{cases} A = diag(1, 2, \cdots, n) + rL^T, \\ B = 2^{-t}I + diag(1, 2, \cdots, n) + rL^T + 2^{-t}L, \end{cases} \tag{4.1}$$

*with L the strictly lower triangular matrix having ones in the lower triangle part and t is a problem parameter to be specified in actual computations.*

The Sylvester equation in Example 4.2 is solved by the ADI, the IADI and the IHSS iterative methods, respectively. Parameters $\alpha_{\exp}^1$ and $\beta_{\exp}^1$ are the numerically found optimal values for IHSS iteration. The corresponding results are listed in Table 4.4, where we set $t = n, r = \frac{1}{n}$, and we use the GMRES as the inner iteration scheme and set $\varepsilon_k = \eta_k = 0.01$, for $k = 0, 1, 2, \cdots$.

From Table 4.4, we can observe that the IADI iteration is usually superior to the ADI iteration in terms of computation efficiency. Moreover, the IADI iteration as well as the ADI iteration outperform the IHSS iteration in terms of both the number of iterations and the computation efficiency, especially when the order of the coefficient matrices is large enough.

Table 4.4 IT and CPU for ADI, IADI and IHSS

| $n$ | ADI | | | | IADI | | | | IHSS | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\alpha_{\exp}$ | $\beta_{\exp}$ | IT | CPU | $\alpha_{\exp}$ | $\beta_{\exp}$ | IT | CPU | $\alpha_{\exp}^1$ | $\beta_{\exp}^1$ | IT | CPU |
| 8 | 3.7 | 1.9 | 9 | 0.010 | 3.7 | 2.1 | 9 | 0.060 | 5 | 1 | 16 | 0.003 |
| 16 | 5.0 | 3.5 | 12 | 0.010 | 5.0 | 3.5 | 12 | 0.060 | 8 | 1 | 21 | 0.028 |
| 32 | 6.7 | 6.1 | 16 | 0.020 | 6.7 | 6.1 | 16 | 0.090 | 13 | 1 | 27 | 0.130 |
| 64 | 9.0 | 8.7 | 23 | 0.060 | 9.0 | 8.7 | 23 | 0.120 | 20 | 3 | 36 | 0.419 |
| 128 | 12.3 | 12.3 | 33 | 0.510 | 12.3 | 12.3 | 33 | 0.731 | 30 | 7 | 45 | 2.321 |
| 256 | 17.0 | 16.9 | 48 | 3.113 | 17.0 | 16.9 | 48 | 3.531 | 50 | 10 | 54 | 19.224 |
| 512 | 23.6 | 23.5 | 69 | 311.325 | 23.6 | 23.5 | 69 | 27.188 | 65 | 11 | 91 | 193.825 |

**5. Conclusions.** In this paper, we have revisited ADI iteration for solving the continuous Sylvester equation $AX + XB = C$, where the coefficient matrices $A$ and $B$ are assumed to be positive semi-definite matrices (not necessarily Hermitian), and at least one of them to be positive definite. For such a class of Sylvester equations, we have analyzed the convergence of the ADI iteration and derived an upper bound of its contraction factor. To reduce the computational complexity, we have also proposed the IADI iteration whose convergence has been analyzed in detail. We have presented some numerical experiments to illustrate the effectiveness of both ADI and IADI iterations.

In order to get the practical choices of the iteration parameters $\alpha$ and $\beta$, we first get the maximal and minimal eigenvalues of $A$ and $B$ via power method and inverse iteration, then get $\tau$ and $\Delta$ according to formulas (2.6) to (2.10), and finally get the approximations to $\alpha$ and $\beta$. The scalars $\alpha_{\exp}$ and $\beta_{\exp}$ in text are obtained by searching the optimal values of the iteration parameters for the ADI iteration in two intervals $(\alpha - 1, \alpha + 1)$ and $(\beta - 1, \beta + 1)$ with stepsize 0.1, respectively. We notice that those iteration parameters $\alpha_{\exp}$ and $\beta_{\exp}$ used in our numerical experiments for the ADI iteration are the experimentally locally optimal values, not the globally optimal values. Therefore, we can expect the ADI and IADI iterations with exact optimal iteration parameters to have better convergence behavior and computational efficiency. However, it is an important and hard task to find the optimal $\alpha$ and $\beta$ which strongly depend on the specific structures and properties of the coefficient matrices $A$ and $B$ and need further in-depth study from the viewpoint of both theory and computations.

We remark here that our convergence theory regarding ADI iterations with two parameters for solving Sylvester equations can be easily extended to solve the general positive definite linear system, $Kx = b$, if $K = U + V$ with $U$ and $V$ being positive definite. In this sense, we can say we have generalized the convergence theory in [15] concerning ADI iterations with two parameters for solving Hermitian positive definite linear system $Kx = b$ to the case of $K$ being non-Hermitian positive definite.

## REFERENCES

[1] Bai ZZ, On Hermitian and skew-Hermitian splitting iteration methods for continuous Sylvester equations, J. Comput Math., 2011; 29: 185-198.

[2] Bai ZZ, A class of modified block SSOR preconditioners for symmetric positive definite systems of linear equations, Adv. Comput. Math., 1999; 10: 169-186.

[3] Bai ZZ, Modified block SSOR preconditioners for symmetric positive definite linear systems, Ann. Operations Research, 2001; 103: 263-282.

[4] Bai ZZ, On SSOR-like preconditioners for non-Hermitian positive definite matrices, Numer. Linear Algebra Appl., 2016; 23: 37-60.

[5] Bai ZZ, Golub GH, Ng MK, Hermitian and skew-Hermitian splitting methods for non-Hermitian positive definite linear systems, SIAM J. Matrix Anal. Appl., 2003; 24: 603-626.

[6] Bai ZZ, Golub GH, Ng MK, On successive-overrelaxation acceleration of the Hermitian and skew-Hermitian splitting iterations, Numer. Linear Algebra Appl., 2007; 14: 319-335.

[7] Bai ZZ, Ng MK, Erratum, Numer. Linear Algebra Appl., 2012; 19: 891.

[8] Bai ZZ, Golub GH, Lu LZ, Yin JF, Block triangular and skew-Hermitian splitting methods for positive-definite linear systems, SIAM J. Sci. Comput., 2005; 26: 844-863.

[9] Bai ZZ, Guo XX, Xu SF, Alternately linearized implicit iteration methods for the minimal nonnegative solutions of the nonsymmetric algebraic Riccati equations, Numer. Linear Algebra Appl., 2006; 13: 655-674.

[10] Bai ZZ, Yin JF, Su YF A shift-splitting preconditioner for non-Hermitian positive definite matrices, J. Comput. Math., 2006; 24: 539-552.

[11] Bai ZZ, Rozloznik M, On the numerical behavior of matrix splitting iteration methods for solving linear systems, SIAM J. Numer. Anal., 2015; 53: 1716-1737.

[12] Bartels RH, Stewart GW, Solution of the matrix equation $AX + XB = C$, Commun. ACM., 1972; 15: 820-826.

[13] Benner P, Li RC, Truhar N, On the ADI method for Sylvester equations, J. Comput. Appl. Math., 2009; 233: 1035-1045.

[14] Bhatia R, Rosenthal P, How and why to solve the operator equation $AX - XB = Y$, Bull. Lond. Math. Soc., 1997; 29: 1-21.

[15] Birkhoff G, Varga RS, Young DM, Alternating direction implicit methods, Adv. Comput., 1962; 3: 189-273.

[16] Golub GH, Nash S, Van Loan CF, A Hessenberg-Schur Method for the problem $AX + XB = C$, IEEE Trans. Automat. Contr., 1979; 24: 909-913.

[17] Golub GH, Van Loan CF, Matrix Computations, 3rd edition, Johns Hopkins University Press, Baltimore, Maryland, 1996.

[18] Gu CQ, Xue HY, A shift-splitting hierarchical identification method for solving Lyapunov matrix equations, Linear

Algebra Appl., 2009; 430: 1517-1530.

[19] Horn RA, Johnson CR, Topics in Matrix Analysis, Cambridge University Press, Cambridge, UK, 1991.

[20] Kürschner P, Benner P, Saak J, Self-generating and efficient shift parameters in ADI methods for large Lyapunov and Sylvester equations, Electr. Trans. Numer. Anal., 2014; 43: 142-162.

[21] P. Lancaster and M. Tismenetsky, The Theory of Matrices, 2nd edition, Academic Press, Orlando, 1985.

[22] N. Levenberg and L. Reichel, A generalized ADI iterative method, Numer. Math., 1993; 66: 215-233.

[23] Liu ZY, Wu NC, Qin XR, Zhang YL, Trigonometric transform splitting methods for real symmetric Toeplitz systems, Comput. Math. Appl., 2018; 75: 2782-2794.

[24] Lu A, Wachspress EL, Solution of Lyapunov equations by alternating direction implicit iteration, Comput. Math. Appl., 1991; 21: 43-58.

[25] Wachspress EL, Trail to a Lyapunov equation solver, Comput. Math. Appl., 2008; 55: 1653-1659.

[26] Saad Y, Iterative Methods for Sparse Linear Systems, 2nd edition, SIAM, Philadelphia, PA, USA, 2003.

[27] Smith RA, Matrix equation $XA + BX = C$, SIAM J. Appl. Math., 1968; 16: 198-201.

[28] Wang X, Li WW, Mao LZ, On positive-definite and skew-Hermitian splitting iterative methods for continuous Sylvester equation $AX + XB = C$, Comput. Math. Appl., 2013; 66: 2352-2361.

[29] Zheng QQ, Ma CF, On normal and skew-Hermitian splitting iterative methods for large sparse continuous Sylvester equations, J. Comput. Appl. Math., 2014; 268: 145-154.