Peter Benner      Grece El Khoury      Miloud Sadkane

# On the Squared Smith Method for Large-Scale Stein Equations

# Max Planck Institute Magdeburg
# Preprints

# On the Squared Smith Method for Large-Scale Stein Equations

Peter Benner*    Grece El Khoury†    Miloud Sadkane‡

September 11, 2012

### Abstract

A squared Smith type algorithm for solving large-scale discrete-time Stein equations is developed. The algorithm uses restarted Krylov spaces to compute approximations of the squared Smith iterations in low-rank factored form. Fast convergence results when very few iterations of the alternating direction implicit method are applied to the Stein equation beforehand. The convergence of the algorithm is discussed and its performance is demonstrated by several test examples.

**Key words.** Stein equation; squared Smith iteration; block-Arnoldi; low-rank factor; ADI iteration.

## 1 Introduction

The Stein equation

$$X - AXB^T = C, \tag{1}$$

where $A$, $B$ and $C$ are given real matrices and $X$ is the unknown plays an important role in areas such as discrete-time control, model reduction of discrete-time dynamical systems, and restoration of images, see, e.g. [12, 14, 2, 9, 5]. It has a unique solution if and only if $\lambda\mu \neq 1$ for all $(\lambda, \mu) \in \Lambda(A) \times \Lambda(B)$, where $\Lambda(S)$ denotes the set of eigenvalues of the square matrix $S$. When the size of $B$ is small compared to that of $A$, an efficient algorithm for solving (1) is based on the alternating direction implicit method (ADI) [5]. Note that when $B$ is small and well conditioned, equation (1) can

---

*Max Planck Institute for Dynamics of Complex Technical Systems, Magdeburg, Germany. `benner@mpi-magdeburg.mpg.de`.

†Université de Brest, CNRS – UMR 6205, Laboratoire de Mathématiques de Bretagne Atlantique. 6, Av. Le Gorgeu. 29238 Brest Cedex 3. France.

‡Université de Brest, CNRS - UMR 6205, Laboratoire de Mathématiques de Bretagne Atlantique. 6, Av. Le Gorgeu. 29238 Brest Cedex 3. France. `miloud.sadkane@univ-brest.fr`.

be reduced to the classical Sylvester equation $AX - X\left(B^T\right)^{-1} = -C\left(B^T\right)^{-1}$ for which efficient methods already exist, see, e.g. [11, 20, 8, 16].

In the present paper we are interested in the case when $A$ and $B$ are large $n \times n$ matrices and $C$ has the low-rank form $C = EF^T$, where $E$ and $F$ are $n \times p$ matrices with $p = \mathrm{rank}(E) = \mathrm{rank}(F) \ll n$. We assume that $A$ and $B$ are discrete-stable (their eigenvalues lie inside the open unit disk). Under this assumption, the unique solution of (1) is given by

$$X = \sum_{j=0}^{\infty} A^j E F^T (B^j)^T. \tag{2}$$

In fact, the existence of (2) is ensured under the weaker condition that the spectral radii $\rho(A)$ and $\rho(B)$ of $A$ and $B$ satisfy $\rho(A)\rho(B) < 1$. Note that if $\rho(A) < 1 < \rho(B) < 1/\rho(A)$, then we can find $\xi$ such that $\rho(B) < \xi < 1/\rho(A)$ and replace $A$ and $B$ by $\xi A$ and $B^T/\xi$ in equation (1).

We assume further that the norms of $A^j$ and $B^j$ decrease as $j$ increases, which ensures that the solution can be approximated as $X \approx Z_1 Z_2^T$ where $Z_1$ and $Z_2$ have ranks much smaller than $n$, see Section 2.

The main aim of the present paper is to describe an effective squared Smith type algorithm for computing the factors $Z_1$ and $Z_2$. The squared Smith method is an improvement of Smith's method initially devised for solving Sylevester equations [21]. In its standard from, the squared Smith method converges quadratically, but necessitates, at each iteration, $n \times n$ matrix-matrix operations, which is the reason why it is mainly confined to matrix equations of small sizes. In this paper we show that these difficulties can be bypassed. We shall use a Krylov space with a restarting scheme to generate the squared Smith iterations in low rank factors. This leads to a computational cost which is linear in $n$, but the quadratic convergence of the original squared Smith method may not be maintained. However, the convergence can be made fast when a simple version of the ADI method, similar to the one proposed in [4], is applied beforehand. In the special case where $B = A$ and $E = F$, an adaptation of squared Smith is developed in [19]. We show that many ideas in this paper can be adapted here but the presence of $B \neq A$ and $E \neq F$ requires a careful implementation.

The paper is organized as follows. Section 2 gives sufficient conditions under which a low-rank solution may be expected and derives low-rank squared Smith iterations from a combination of the squared Smith method and a block variant of the Arnoldi algorithm. Section 3 estimates the error between the exact solution and its low rank approximation and expresses the residual in terms of quantities readily computable. The residual is used in Section 4 to develop a cheap restarting scheme to overcome the increase in computational cost and memory requirements for the block Arnoldi bases. The acceleration of convergence due to a simple version of ADI iterations is discussed in Section 5. Section 6 is devoted to numerical tests, and concluding remarks are given in Section 7.

## 2 Low-rank approximation

The solution (2) can be written as

$$X = \Big( E, AE, A^2E, \dots \Big) \Big( F, BF, B^2F, \dots, \Big)^T,$$

and due to the rank properties and the Cayley-Hamilton theorem, we have

$$\mathrm{rank}(X) \leq \min \Big( \mathrm{rank}\left( E, AE, \dots, A^{n-1}E \right), \mathrm{rank}\left( F, AF, \dots, B^{n-1}F \right) \Big). \quad (3)$$

The ranks in the right-hand sides of (3) can be much smaller than $n$, depending on the proprieties of $A$, $B$, $E$, and $F$. For example, if the columns of $E$ (or $F$) span an invariant subspace of $A$ (or $B$), then $\mathrm{rank}(X) \leq p$. If the norms of powers of $A$ (or $B$) decrease rapidly, then $\mathrm{rank}(X)$ is small compared to $n$. In the general case, $X$ can be decomposed as

$$X = X_{k,1} + X_{k,2},$$

where $X_{k,1} = \sum_{j=0}^{k-1} A^j E F^T (B^j)^T$ and $X_{k,2} = X - X_{k,1}$ and $k$ is chosen such that $kp \leq n - 1$. Since

$$X_{k,1} = \Big( E, AE, \dots, A^{k-1}E \Big) \Big( F, AF, \dots, A^{k-1}F \Big)^T,$$

it is clear that $\mathrm{rank}(X_{1,k}) \leq pk$. Let the singular values of $X$ be labeled so that $\sigma_1(X) \geq \sigma_2(X) \dots \geq \sigma_n(X)$. Then the Schmidt-Mirsky theorem [22] gives $\sigma_{kp+1}(X) \leq \|X - X_{k,1}\|$, and therefore

$$\begin{aligned} \sigma_{kp+1}(X) &\leq \left\| \sum_{j=k}^{\infty} A^j E F^T (B^j)^T \right\| \\ &= \left\| A^k X (B^k)^T \right\| \leq \|A^k\| \|B^k\| \|X\|. \end{aligned}$$

This yields the upper bound

$$\frac{\sigma_{kp+1}(X)}{\sigma_1(X)} \leq \|A^k\| \|B^k\|. \quad (4)$$

None of the bounds (3), (4) is sharp but they both show that a solution with a (numerical) low-rank can be expected provided that the norms of the powers of $A$ or $B$ decrease rapidly.

### 2.1 Low-rank squared Smith approximation

The solution $X$ can be approximated by the partial sum

$$X_k = \sum_{j=0}^{2^k-1} A^j E F^T (B^T)^j \quad (5)$$

3

with large $k$.

We see that $X_0 = EF^T$ and for $k \geq 1$

$$
\begin{aligned}
X_k &= \sum_{j=0}^{2^{k-1}-1} A^j E F^T (B^T)^j + A^{2^{k-1}} \left( \sum_{j=0}^{2^{k-1}-1} A^j E F^T (B^T)^j \right) (B^T)^{2^{k-1}}, \\
&= X_{k-1} + A^{2^{k-1}} X_{k-1} (B^T)^{2^{k-1}}
\end{aligned}
$$

and

$$
A^{2^k} = \left( A^{2^{k-1}} \right)^2, \quad (B^T)^{2^k} = \left( (B^T)^{2^{k-1}} \right)^2.
$$

Hence, the squared Smith scheme can be written as

$$
X_0 = EF^T, \ A_0 = A, \ B_0 = B,
$$

$$
X_k = X_{k-1} + A_{k-1} X_{k-1} B_{k-1}^T, \ A_k = A_{k-1}^2, \ B_k = B_{k-1}^2, \ k \geq 1. \tag{6}
$$

It is clear that (6) should not be used as such since the matrix sequences $A_k$ and $B_k$ are large and dense. However, from (5) we have

$$
X_k = \left( E, AE, \dots, A^{2^k-1}E \right) \left( F, BF, \dots, B^{2^k-1}F \right)^T \tag{7}
$$

and then

$$
X_k \approx Z_k^E (Z_k^F)^T, \tag{8}
$$

where $Z_k^E$ and $Z_k^F$ are matrices of small rank that can be constructed from the Krylov spaces

$$
\begin{aligned}
\mathcal{K}_k(A, E) &= \text{range} \left( E, AE, \dots, A^{2^k-1}E \right), \\
\mathcal{K}_k(A, F) &= \text{range} \left( F, BF, \dots, B^{2^k-1}F \right).
\end{aligned}
$$

A natural way to compute $Z_k^E$ and $Z_k^F$ is the block Arnoldi algorithm applied to $A$ and $B$ and starting with $E$ and $F$, respectively. The following algorithm is the version applied to $A$ and started with $E$.

For the implementation issues, see the discussion after Algorithm 2.

For $j = 1, \dots, 2^m$, let

$$
\mathbb{Q}_j^E = (Q_1^E, \ \dots, \ Q_j^E), \quad \mathbb{Q}_{j+1}^E = (\mathbb{Q}_j^E \quad Q_{j+1}^E),
$$

$$
\mathbb{H}_j^E = (H_{i,l}^E)_{1 \leq i, l \leq j}, \quad \underline{\mathbb{H}}_j^E = \begin{pmatrix} \mathbb{H}_j^E \\ H_{j+1,j}^E I_j^T \end{pmatrix},
$$

where $I_j^T = (0\,0\,\dots\,I)$ and $0$ and $I$ denote the zero and identity matrices of appropriate sizes.

---
**Algorithm 1** Block Arnoldi method
---
INPUT: $A \in \mathbb{R}^{n,n}, E \in \mathbb{R}^{n,p}$, an integer $m$ such that $2^m \ll n$.
OUTPUT: Arnoldi basis $(Q_1^E, \ldots, Q_{2^m+1}^E)$, and blocks $H_{i,j}^E$ of corresponding Hessenberg matrix.

1:  $QR$ factorize $E = Q_1^E R_1^E$
2:  **for** $j = 1, \ldots, 2^m$ **do**
3:      $W_j^E = AQ_j^E$
4:      **for** $i = 1, \ldots, j$ **do**
5:          $H_{i,j}^E = (Q_i^E)^T W_j^E$
6:          $W_j^E := W_j^E - Q_i^E H_{i,j}^E$
7:      **end for** i
8:      $QR$ factorize $W_j^E = Q_{j+1}^E H_{j+1,j}^E$
9:  **end for** j
---

It is known (see, e.g. [17]) that $\mathbb{H}_j^E$ is block upper Hessenberg and the columns of $\mathbb{Q}_{2^m}^E$ form an orthonormal basis of $\mathbb{K}_m(A, Q_1^E)$ if no deflation occurs (which we assume for clarity of presentation — in the actual implementation, this is taken care of). From Algorithm 1 and the equalities above we obtain

$$A\mathbb{Q}_j^E = \mathbb{Q}_{j+1}^E \underline{\mathbb{H}}_j^E. \tag{9}$$

We will often need operations of type $A_l \mathbb{Q}_j^E$ where $A_l = A_{l-1}^2 = A^{2^l}$ is defined in (6). From (9) we have

$$A_l \mathbb{Q}_j^E = \underbrace{A\big(A\big(\ldots\big(A\,\mathbb{Q}_j^E\big)\big)\big)}_{2^l \text{ times}} = \mathbb{Q}_{j+2^l}^E \prod_{k=j+2^l-1}^{j} \underline{\mathbb{H}}_k^E. \tag{10}$$

The same algorithm applied to $B$ and started with $F$ leads to similar formulas as above. We will refer to these by the superscript $F$.

The approximation (8) is then obtained as follows:

$$X_0 = EF^T \equiv Z_0^E Z_0^{F^T} \text{ with } Z_0^E = Q_1^E R_1^E \text{ and } Z_0^F = Q_1^F R_1^F, \tag{11}$$

$$X_1 = X_0 + AX_0 B^T = (E, AE)(F, BF)^T \tag{12}$$

with

$$
\begin{aligned}
(E, AE) &= (Q_1^E R_1^E, AQ_1^E R_1^E) \\
&= (Q_1^E R_1^E, \mathbb{Q}_2^E \underline{\mathbb{H}}_1^E R_1^E) \\
&= \mathbb{Q}_2^E \left( \begin{pmatrix} R_1^E \\ 0 \end{pmatrix}, \underline{\mathbb{H}}_1^E R_1^E \right).
\end{aligned}
$$

Similarly, for $(F, BF)$ we have

$$(F, BF) = \mathbb{Q}_2^F \left( \begin{pmatrix} R_1^F \\ 0 \end{pmatrix}, \mathbb{H}_1^F R_1^F \right).$$

The low-rank approximation of $(E, AE)$ and $(F, AF)$ are then obtained through SVDs of the small matrices [1] $\left( \begin{pmatrix} R_1^E \\ 0 \end{pmatrix}, \mathbb{H}_1^E R_1^E \right)$ and $\left( \begin{pmatrix} R_1^F \\ 0 \end{pmatrix}, \mathbb{H}_1^F R_1^F \right)$ obtained by deleting the singular values which are smaller than a threshold $\text{tol}_1$:

$$\left( \begin{pmatrix} R_1^E \\ 0 \end{pmatrix}, \mathbb{H}_1^E R_1^E \right) = U_1^E S_1^E (V_1^E)^T + \Delta_1^E, \tag{13}$$

$$\left( \begin{pmatrix} R_1^F \\ 0 \end{pmatrix}, \mathbb{H}_1^F R_1^F \right) = U_1^F S_1^F (V_1^F)^T + \Delta_1^F, \tag{14}$$

with

$$\|\Delta_1^E\| < \text{tol}_1 \quad \text{and} \quad \|\Delta_1^F\| < \text{tol}_1, \tag{15}$$

where the matrices $U_1^E$, $U_1^F$, $V_1^E$ and $V_1^F$ have orthonormal columns, and $S_1^E$ and $S_1^F$ are diagonal matrices whose diagonals contain the singular values which are larger than $\text{tol}_1$. Here and throughout the paper, the symbol $\| \ \|$ denotes the spectral norm.

Care must be taken when deleting the smaller singular values to make the operation $(V_1^E)^T V_1^F$ possible. In our implementation, the number of singular values deleted equals the maximum of the number of singular values that are smaller than $\text{tol}_1$ in both $\left( \begin{pmatrix} R_1^E \\ 0 \end{pmatrix}, \mathbb{H}_1^E R_1^E \right)$ and $\left( \begin{pmatrix} R_1^F \\ 0 \end{pmatrix}, \mathbb{H}_1^F R_1^F \right)$. Then we have the first low-rank approximation:

$$X_1 \approx Z_1^E (Z_1^F)^T \quad \text{with} \quad Z_1^E = \mathbb{Q}_2^E U_1^E S_1^E (V_1^E)^T V_1^F \quad \text{and} \quad Z_1^F = \mathbb{Q}_2^F U_1^F S_1^F. \tag{16}$$

Note that the choice $Z_1^E = \mathbb{Q}_2^E U_1^E S_1^E (V_1^E)^T$ and $Z_1^F = \mathbb{Q}_2^F U_1^F S_1^F (V_1^F)^T$ is not recommended since the numbers of columns of the factors $Z_k^E$ and $Z_k^F$ are intended to increase at each iteration $k$, see Algorithm 2.

The same procedure leads us to the second step:

$$\begin{aligned}
X_2 &= X_1 + A_1 X_1 B_1^T \\
&\approx Z_1^E Z_1^{F^T} + A_1 Z_1^E Z_1^{F^T} B_1^T \\
&= (Z_1^E, A_1 Z_1^E)(Z_1^F, B_1 Z_1^F)^T
\end{aligned}$$

with

$$\begin{aligned}
(Z_1^E, A_1 Z_1^E) &= (\mathbb{Q}_2^E U_1^E S_1^E (V_1^E)^T V_1^F, A^2 \mathbb{Q}_2^E U_1^E S_1^E (V_1^E)^T V_1^F) \\
&= (\mathbb{Q}_2^E U_1^E S_1^E (V_1^E)^T V_1^F, \mathbb{Q}_4^E \mathbb{H}_3^E \mathbb{H}_2^E U_1^E S_1^E (V_1^E)^T V_1^F) \\
&= \mathbb{Q}_4^E \left( \begin{pmatrix} U_1^E S_1^E (V_1^E)^T V_1^F \\ 0 \end{pmatrix}, \mathbb{H}_3^E \mathbb{H}_2^E U_1^E S_1^E (V_1^E)^T V_1^F \right).
\end{aligned}$$

---

[1] A cheaper alternative is the rank revealing QR algorihtm, see e.g. [1]

6

and similarly for $(Z_1^F, B_1 Z_1^F)$:

$$(Z_1^F, B_1 Z_1^F) = \mathbb{Q}_4^F \left( \begin{pmatrix} U_1^F S_1^F \\ 0 \end{pmatrix}, \mathbb{H}_3^F \mathbb{H}_2^F U_1^F S_1^F \right).$$

As before the SVDs

$$\left( \begin{pmatrix} U_1^E S_1^E (V_1^E)^T V_1^F \\ 0 \end{pmatrix}, \mathbb{H}_3^E \mathbb{H}_2^E U_1^E S_1^E (V_1^E)^T V_1^F \right) = U_2^E S_2^E (V_2^E)^T + \Delta_2^E, \qquad (17)$$

$$\left( \begin{pmatrix} U_1^F S_1^F \\ 0 \end{pmatrix}, \mathbb{H}_3^F \mathbb{H}_2^F U_1^F S_1^F \right) = U_2^F S_2^F (V_2^F)^T + \Delta_2^F, \qquad (18)$$

with a threshold $\mathrm{tol}_2$ such that

$$\|\Delta_2^E\| < \mathrm{tol}_2 \quad \text{and} \quad \|\Delta_2^F\| < \mathrm{tol}_2 \qquad (19)$$

lead to the second low-rank approximation:

$$X_2 \approx Z_2^E Z_2^{F^T} \text{ with } Z_2^E = \mathbb{Q}_4^E U_2^E S_2^E (V_2^E)^T V_2^F \text{ and } Z_2^F = \mathbb{Q}_4^F U_2^F S_2^F. \qquad (20)$$

More generally, at step $k$ we have:

$$X_k \approx (Z_{k-1}^E, A_{k-1} Z_{k-1}^E)(Z_{k-1}^F, B_{k-1} Z_{k-1}^F)^T$$

with

$$(Z_{k-1}^E, A_{k-1} Z_{k-1}^E) =$$

$$(\mathbb{Q}_{2^k-1}^E U_{k-1}^E S_{k-1}^E (V_{k-1}^E)^T V_{k-1}^F, A^{2^{k-1}} \mathbb{Q}_{2^k-1}^E U_{k-1}^E S_{k-1}^E (V_{k-1}^E)^T V_{k-1}^F) =$$

$$\mathbb{Q}_{2^k}^E \left( \begin{pmatrix} U_{k-1}^E S_{k-1}^E (V_{k-1}^E)^T V_{k-1}^F \\ 0 \end{pmatrix}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^E U_{k-1}^E S_{k-1}^E (V_{k-1}^E)^T V_{k-1}^F \right)$$

and

$$\begin{aligned} (Z_{k-1}^F, B_{k-1} Z_{k-1}^F) &= (\mathbb{Q}_{2^k-1}^F U_{k-1}^F S_{k-1}^F, B^{2^{k-1}} \mathbb{Q}_{2^k-1}^F U_{k-1}^F S_{k-1}^F) \\ &= \mathbb{Q}_{2^k}^F \left( \begin{pmatrix} U_{k-1}^F S_{k-1}^F \\ 0 \end{pmatrix}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^F U_{k-1}^F S_{k-1}^F \right). \end{aligned}$$

Then we compute, with a tolerance $\mathrm{tol_k}$, the reduced SVDs

$$\left( \begin{pmatrix} U_{k-1}^E S_{k-1}^E (V_{k-1}^E)^T V_{k-1}^F \\ 0 \end{pmatrix}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^E U_{k-1}^E S_{k-1}^E (V_{k-1}^E)^T V_{k-1}^F \right) =$$

$$U_k^E S_k^E (V_k^E)^T + \Delta_k^E, \qquad (21)$$

7

$$\left( \begin{pmatrix} U_{k-1}^F S_{k-1}^F \\ 0 \end{pmatrix}, \prod_{j=2^k-1}^{2^{k-1}} \underline{\mathbb{H}}_j^F U_{k-1}^F S_{k-1}^F \right) = U_k^F S_k^F (V_k^F)^T + \Delta_k^F \tag{22}$$

with

$$\|\Delta_k^E\| < \mathrm{tol_k} \quad \text{and} \quad \|\Delta_k^F\| < \mathrm{tol_k} \tag{23}$$

and we obtain the low rank approximation

$$X_k \approx Z_k^E (Z_k^F)^T \text{ with } Z_k^E = \mathbb{Q}_{2^k}^E U_k^E S_k^E (V_k^E)^T V_k^F \text{and } Z_k^F = \mathbb{Q}_{2^k}^F U_k^F S_k^F. \tag{24}$$

# 3 Error and residual estimates

In this section we estimate the error

$$X - Z_k^E (Z_k^F)^T \tag{25}$$

between the exact and the computed solution and the residual

$$\Gamma_k = EF^T + A Z_k^E (Z_k^F)^T B^T - Z_k^E (Z_k^F)^T \tag{26}$$

associated with the computed solution.

The estimates in Propositions 1–4, though complicated to compute, will provide a reasonable indication of the error due to the Krylov space approximation and the SVD truncations. Proposition 4 shows that the norm of the residual can be computed at a lower cost. As we will see in Algorithm 2, it will be used to stop and/or restart the iterations.

## 3.1 Estimation of the error

From (2) and (5) we obtain

$$X - X_k = \sum_{j \geq 2^k} A^j EF^F (B^j)^T. \tag{27}$$

Since the eigenvalues of $A$ lie in the open unit disk, they are certainly inside some circle $|z| = \rho_a$ of center 0 and radius $\rho_a < 1$. The Cauchy integral formula [10] gives

$$A^j = \frac{1}{2i\pi} \int_{|z|=\rho_a} z^j (zI - A)^{-1} dz = \frac{\rho_a^{j+1}}{2\pi} \int_0^{2\pi} e^{i(j+1)\theta} (\rho_a e^{i\theta} I - A)^{-1} d\theta$$

and a similar formula for $B^j$ can be derived. Thus

$$\|A^j\| \leq C_a \rho_a^{j+1}, \quad \|B^j\| \leq C_b \rho_b^{j+1} \tag{28}$$

with

$$C_a = \max_{0 \leq \theta \leq 2\pi} \|(\rho_a e^{i\theta} I - A)^{-1}\|, \quad C_b = \max_{0 \leq \theta \leq 2\pi} \|(\rho_b e^{i\theta} I - B)^{-1}\|.$$

Taking the norm in (27) and using (28) we obtain the following proposition.

**Proposition 1.** *For $k \geq 0$ we have*

$$\|X - X_k\| \leq \frac{C_a C_b \|EF^T\|}{1 - \rho_a \rho_b} \left(\rho_a \rho_b\right)^{2^k + 1}.$$

This proposition holds true if we only assume that $\rho_a \rho_b < 1$ and it clearly shows that $\|X - X_k\|$ tends to zero as $k$ increases. However, the convergence may be slowed down if $\rho_a \rho_b$ is very close to 1 or if the constants $C_a$ or $C_b$ are too large. In the first case, a simple variant of the ADI iteration can be used to minimize the spectral radii, see Section 5. In the second case, the pseudospectra of $A$ or $B$ may significantly protrude from the unit circle [23] and a low-rank approximate solution may not exist.

Next, we estimate the error between $X_k$ and $Z_k^E (Z_k^F)^T$.

**Proposition 2.** *For $k \geq 0$ we have*

$$\|X_k - Z_k^E (Z_k^F)^T\| \leq \mu_k$$

*with $\mu_0 = 0$ and for $k \geq 1$*

$$\mu_k \leq \delta_k + \left\| \left( \begin{pmatrix} I \\ 0 \end{pmatrix}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^E \right) \right\| \left\| \left( \begin{pmatrix} I \\ 0 \end{pmatrix}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^F \right) \right\| \delta_{k-1} \tag{29}$$

*with*

$$\delta_0 = 0 \quad \text{and} \quad \text{for} \quad k \geq 1, \quad \delta_k = \left( \|S_k^E\| + \|S_k^F\| \right) \text{tol}_k + \text{tol}_k^2,$$

*where $\text{tol}_1, \text{tol}_2, \ldots$ are the SVD thresholds defined in (15), (19), ..., (23).*

*Proof.* For $k = 0$, the bound is satisfied since $X_0 = EF^T = Z_0^E (Z_0^F)^T$.

For $k = 1$, we have from (13) and (14)

$$
\begin{aligned}
(E, AE) &= \mathbb{Q}_2^E \left( U_1^E S_1^E (V_1^E)^T + \Delta_1^E \right), \\
(F, BF) &= \mathbb{Q}_2^F \left( U_1^F S_1^F (V_1^F)^T + \Delta_1^F \right).
\end{aligned}
$$

Hence

$$
\begin{aligned}
X_1 &= (E, AE)(F, BF)^T \\
&= \mathbb{Q}_2^E \left( U_1^E S_1^E (V_1^E)^T + \Delta_1^E \right) \left( U_1^F S_1^F (V_1^F)^T + \Delta_1^F \right)^T (\mathbb{Q}_2^F)^T \\
&= Z_1^E (Z_1^F)^T + \widetilde{X}_1
\end{aligned}
$$

with

$$\widetilde{X}_1 = \mathbb{Q}_2^E \left( U_1^E S_1^E (V_1^E)^T (\Delta_1^F)^T + \Delta_1^E V_1^F S_1^F (U_1^F)^T + \Delta_1^E (\Delta_1^F)^T \right) (\mathbb{Q}_2^F)^T$$

and hence

$$\|\widetilde{X}_1\| \leq \left( \|S_1^E\| + \|S_1^F\| \right) \text{tol}_1 + \text{tol}_1^2 = \delta_1,$$

which shows the proposition for $k = 1$.

The general case is straightforward but tedious. From (21)–(23) we have

$$\left( E, AE, \dots, A^{2^{k-1}} E \right) =$$

$$\mathbb{Q}_{2^k}^E \left( \begin{pmatrix} U_{k-1}^E S_{k-1}^E (V_{k-1}^E)^T \\ 0 \end{pmatrix}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^E U_{k-1}^E S_{k-1}^E (V_{k-1}^E)^T \right) +$$

$$\mathbb{Q}_{2^k}^E \left( \begin{pmatrix} \Delta_{k-1}^E \\ 0 \end{pmatrix}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^E \Delta_{k-1}^E \right)$$

and

$$\left( F, BF, \dots, B^{2^{k-1}} F \right) =$$

$$\mathbb{Q}_{2^k}^F \left( \begin{pmatrix} U_{k-1}^F S_{k-1}^F (V_{k-1}^F)^T \\ 0 \end{pmatrix}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^F U_{k-1}^F S_{k-1}^F (V_{k-1}^F)^T \right) +$$

$$\mathbb{Q}_{2^k}^F \left( \begin{pmatrix} \Delta_{k-1}^F \\ 0 \end{pmatrix}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^F \Delta_{k-1}^F \right).$$

Hence

$$\begin{aligned}
X_k &= \left( E, AE, \dots, A^{2^{k-1}} E \right) \left( F, BF, \dots, B^{(2^{k-1})} F \right)^T \\
&= Z_k^E (Z_k^F)^T + \widetilde{X}_k
\end{aligned}$$

with

$$\widetilde{X}_k = \mathbb{Q}_{2^k}^E \left( \Upsilon_0^{E,F} + \Upsilon_1^E (\Upsilon_2^F)^T + \Upsilon_2^E (\Upsilon_1^F)^T + \Upsilon_2^E (\Upsilon_2^F)^T \right) (\mathbb{Q}_{2^k}^F)^T,$$

where

$$\Upsilon_0^{E,F} = U_k^E S_k^E (V_k^E)^T (\Delta_k^F)^T + \Delta_k^E V_k^F S_k^F (U_k^F)^T + \Delta_k^E (\Delta_k^F)^T$$

and for $G = E$ or $F$

$$\Upsilon_1^G = \left( \begin{pmatrix} I \\ 0 \end{pmatrix}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^G \right) \begin{pmatrix} U_{k-1}^G S_{k-1}^G (V_{k-1}^G)^T & 0 \\ 0 & U_{k-1}^G S_{k-1}^G (V_{k-1}^G)^T \end{pmatrix}$$

$$\Upsilon_2^G = \left( \begin{pmatrix} I \\ 0 \end{pmatrix}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^G \right) \begin{pmatrix} \Delta_{k-1}^G & 0 \\ 0 & \Delta_{k-1}^G \end{pmatrix}.$$

We clearly have

$$\|\Upsilon_0^{E,F}\| \leq \left(\|S_k^E\| + \|S_k^F\|\right)\mathrm{tol_k} + \mathrm{tol_k^2} = \delta_{\mathrm{k}},$$

$$\|\Upsilon_1^G\| \leq \left\|\left(\binom{I}{0}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^G\right)\right\| \|S_{k-1}^G\|,$$

$$\|\Upsilon_2^G\| \leq \left\|\left(\binom{I}{0}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^G\right)\right\| \mathrm{tol_{k-1}},$$

from which the proposition follows. $\qquad\qquad\square$

By combining Propositions 1 and 2 we obtain

**Proposition 3.** *For $k \geq 0$ we have*

$$\|X - Z_k^E\left(Z_k^F\right)^T\| \leq \frac{C_a C_b \|EF^T\|}{1 - \rho_a \rho_b}\left(\rho_a \rho_b\right)^{2^k+1} + \mu_k.$$

When the constant $C_a C_b$ is not large and $\rho_a \rho_b$ is not close to 1, the proposition shows that for large $k$, the error and the residual norms behave respectively as $\mu_k$ and $\|EF^T + AX_k B^T - X_k\|$.

**Remark 3.1.** *In the special case where $B = A$ and $F = E$, the factor $Z_k^E$ is entirely defined from $\left(E, AE, \ldots, A^{2^{k-1}}E\right)$ and this helps to improve the error estimate. To illustrate this point, consider the simple case $k = 1$. Then, from (13)–(15) we have, when $B = A$ and $F = E$*

$$\left(E, AE\right) = \mathbb{Q}_2^E\left(U_1^E S_1^E (V_1^E)^T + \Delta_1^E\right).$$

*Since $\Delta_1^E V_1^E = 0$, we obtain*

$$X_1 = \left(E, AE\right)\left(E, AE\right)^T = Z_1^E(Z_1^E)^T + \mathbb{Q}_2^E \Delta_1^E\left(\mathbb{Q}_2^E \Delta_1^E\right)^T$$

*with $Z_1^E = \mathbb{Q}_2^E U_1^E S_1^E$.*
*In the general case ($B \neq A$, $F \neq E$), the expressions of $Z_1^E$ and $Z_1^F$ are given in (16).*

**Remark 3.2.** *Propositions 1 and 2 produce actually worst-case bounds. In practical computations, the convergence can be much better than these bounds predict.*
*The choice of $\mathrm{tol_k}$ may be difficult to tune, in our experiments we use the same tolerance threshold, that is, for all $k$, $\mathrm{tol_k} = \mathrm{tol_{svd}}$.*

## 3.2 Estimation of the residual

The following proposition relates the residual to the powers of $A$ and $B$ and the error incurred in $X_k - Z_k^E Z_k^F$.

**Proposition 4.** *For $k \geq 0$, we have for the residual $\Gamma_k$ defined in (26),*

$$\|\Gamma_k\| \leq \|E\|\|F\|\left\|A^{2^k}(B^{2^k})^T\right\| + (1 + \|A\|\|B\|)\mu_k,$$

*where $\mu_k$ is defined in Proposition 2.*

*Proof.*

$$\Gamma_k = EF^T + AX_kB^T - X_k + A(Z_k^E(Z_k^F)^T - X_k)B^T - (Z_k^E(Z_k^F)^T - X_k),$$

thus

$$\|\Gamma_k\| \leq \|EF^T + AX_kB^T - X_k\| + (1 + \|A\|\|B\|)\|Z_k^E(Z_k^F)^T - X_k\|.$$

From (5) we obtain

$$EF^T + AX_kB^T - X_k = A^{2^k}EF^T(B^{2^k})^T \tag{30}$$

and the proof is completed by using Proposition 2. $\qquad\square$

The following proposition shows that the norm of the residual involves quantities readily computable. In particular, it will be used as a stopping criterion in Algorithm 2.

**Proposition 5.** *We have $\|\Gamma_0\| = \|(\mathbb{H}_1^E R_1^E)(\mathbb{H}_1^F R_1^F)^T\|$ and for $k \geq 1$,*

$$\|\Gamma_k\| = \left\| \begin{pmatrix} R_1^E \\ 0 \end{pmatrix} \left( (R_1^F)^T \quad 0 \right) + \left( \mathbb{H}_{2^k}^E U_k^E S_k^E (V_k^E)^T V_k^F \right) \left( \mathbb{H}_{2^k}^F U_k^F S_k^F \right)^T - \begin{pmatrix} U_k^E S_k^E (V_k^E)^T V_k^F \\ 0 \end{pmatrix} \left( S_k^F (U_k^F)^T \quad 0 \right) \right\|.$$

*Proof.* We have

$$\Gamma_0 = (AE)(BF)^T = (\mathbb{Q}_2^E \mathbb{H}_1^E R_1^E)(\mathbb{Q}_2^F \mathbb{H}_1^F R_1^F)^T.$$

Hence $\|\Gamma_0\| = \|(\mathbb{H}_1^E R_1^E)(\mathbb{H}_1^F R_1^F)\|$.

For $k \geq 1$, using (26), (24) and (9) the residual $\Gamma_k$ can be written

$$\begin{aligned}
\Gamma_k &= (Q_1^E R_1^E)(Q_1^F R_1^F)^T + A(\mathbb{Q}_{2^k}^E U_k^E S_k^E (V_k^E)^T V_k^F)(\mathbb{Q}_{2^k}^F U_k^F S_k^F)^T B^T - \\
&\quad (\mathbb{Q}_{2^k}^E U_k^E S_k^E (V_k^E)^T V_k^F)(\mathbb{Q}_{2^k}^F U_k^F S_k^F)^T \\
&= \mathbb{Q}_{2^k+1}^E \left[ \begin{pmatrix} R_1^E \\ 0 \end{pmatrix} \left( (R_1^F)^T \quad 0 \right) + \left( \mathbb{H}_{2^k}^E U_k^E S_k^E V_k^{E^T} V_k^F \right) \left( \mathbb{H}_{2^k}^F U_k^F S_k^F \right)^T - \right. \\
&\quad \left. \begin{pmatrix} U_k^E S_k^E V_k^{E^T} V_k^F \\ 0 \end{pmatrix} \left( S_k^F (U_k^F)^T \quad 0 \right) \right] (\mathbb{Q}_{2^k+1}^F)^T
\end{aligned}$$

and the proposition follows by taking the norm. $\qquad\square$

In practice we cannot use Algorithm 1 with large $m$ since the computational cost and memory storage become prohibitive. On the other hand, with a too small $m$, the Krylov spaces thus obtained will not contain sufficient information to allow a good approximation of the solution. To remedy this, we will use in the next two sections a restarting technique and a simple version of the ADI iteration.

# 4 Restarting the low-rank approximation

As in iterative methods for large linear systems [17], restarting is based on the residual. Our residual $\Gamma_k$ has a special form. From (30) we see that its rank is not larger than $p$. Its smallest singular values decrease as the number of iterations increases. We will use this information to construct new $E^{\mathrm{rst}}$ and $F^{\mathrm{rst}}$ that will be used in place of $E$ and $F$ for the next restart. The construction is based on an incomplete SVD of the residual $\Gamma_k$ computed with the help of Proposition 4 at a lower cost.

Consider the reduced SVD

$$\begin{pmatrix} R_1^E \\ 0 \end{pmatrix} \left( \left( R_1^F \right)^T \quad 0 \right) + \left( \mathbb{H}_{2^k}^E U_k^E S_k^E V_k^{E^T} V_k^F \right) \left( \mathbb{H}_{2^k}^F U_k^F S_k^F \right)^T - \\ \begin{pmatrix} U_k^E S_k^E V_k^{E^T} V_k \\ 0 \end{pmatrix} \left( U_k^F S_k^F \quad 0 \right) = U_k S_k V_k^T + \Theta_k,$$

where $U_k$ and $V_k$ have orthonormal columns, $S_k$ is diagonal whose elements are larger than some convergence threshold $\mathrm{tol}_{\mathrm{cvg}}$ and $\Theta_k$ contains the rest of the SVD with $\|\Theta_k\| \le \mathrm{tol}_{\mathrm{cvg}}$. Then from the proof of Proposition 5 we have the decomposition

$$\begin{aligned} \Gamma_k &= \mathbb{Q}_{2^k+1}^E (U_k S_k V_k^T + \Theta_k)(\mathbb{Q}_{2^k+1}^F)^T \\ &= E^{\mathrm{rst}}(F^{\mathrm{rst}})^T + \widetilde{\Theta}_k \end{aligned} \tag{31}$$

with $E^{\mathrm{rst}} = \mathbb{Q}_{2^k+1}^E U_k S_k^{\frac{1}{2}}$, $F^{\mathrm{rst}} = \mathbb{Q}_{2^k+1}^F V_k S_k^{\frac{1}{2}}$ and $\widetilde{\Theta}_k = \mathbb{Q}_{2^k+1}^E \Theta_k (\mathbb{Q}_{2^k+1}^F)^T$, $\|\widetilde{\Theta}_k\| = \|\Theta_k\| \le \mathrm{tol}_{\mathrm{cvg}}$.

Denote

$$Q_1^{\mathrm{rst},E} = \mathbb{Q}_{2^k+1}^E U_k \ , \ Q_1^{\mathrm{rst},F} = \mathbb{Q}_{2^k+1}^F V_k, \ R_1^{\mathrm{rst},E} = R_1^{\mathrm{rst},F} = S_k^{\frac{1}{2}},$$

then we have the $QR$ factorizations

$$E^{\mathrm{rst}} = Q_1^{\mathrm{rst},\,E} R_1^{\mathrm{rst},\,E} \text{ and } F^{\mathrm{rst}} = Q_1^{\mathrm{rst},\,F} R_1^{\mathrm{rst},\,F}. \tag{32}$$

Applying Algorithm 1 with $A$ and $Q_1^{\mathrm{rst},E}$ and with $B$ and $Q_1^{\mathrm{rst},F}$ and proceeding as in Section 2, we obtain factors of the form $Z_k^{\mathrm{rst},E}$ and $Z_k^{\mathrm{rst},F}$. The residual associated with $Z_k^{\mathrm{rst},E} \left( Z_k^{\mathrm{rst},F} \right)^T$ is

$$\Gamma_k^{\mathrm{rst}} = E^{\mathrm{rst}} \left( F^{\mathrm{rst}} \right)^T + A Z_k^{\mathrm{rst},E} \left( Z_k^{\mathrm{rst},F} \right)^T B^T - Z_k^{\mathrm{rst},E} \left( Z_k^{\mathrm{rst},F} \right)^T \tag{33}$$

and its norm is computed by Proposition 5.

If this norm is smaller than $\mathrm{tol}_{\mathrm{cvg}}$, then the iterations can be stopped and the new approximate solution is

$$X_k^{\mathrm{new}} \approx \left( Z_k^E, Z_k^{\mathrm{rst},\,E} \right) \left( Z_k^F, Z_k^{\mathrm{rst},\,F} \right)^T. \tag{34}$$

The corresponding residual is given by

$$\begin{aligned} \Gamma_k^{\mathrm{new}} &= EF + A \left( Z_k^E, Z_k^{\mathrm{rst},\,E} \right) \left( Z_k^F, Z_k^{\mathrm{rst},\,F} \right)^T B^T - \left( Z_k^E, Z_k^{\mathrm{rst},\,E} \right) \left( Z_k^F, Z_k^{\mathrm{rst},\,F} \right)^T \\ &= \Gamma_k + \Gamma_k^{\mathrm{rst}} - E^{\mathrm{rst}} \left( F^{\mathrm{rst}} \right)^T \\ &= \Gamma_k^{\mathrm{rst}} + \widetilde{\Theta}_k, \end{aligned}$$

13

from which we see that $\|\Gamma_k^{\mathrm{new}}\| \le \|\Gamma_k^{\mathrm{rst}}\| + \|\widetilde{\Theta}_k\| \le 2\,\mathrm{tol}_{\mathrm{cvg}}$, so that $\Gamma_k^{\mathrm{new}}$ may have a norm slightly larger than $\mathrm{tol}_{\mathrm{cvg}}$.

If the norm of $\Gamma_k^{\mathrm{rst}}$ is larger than $\mathrm{tol}_{\mathrm{cvg}}$, then $\Gamma_k^{\mathrm{rst}}$ is decomposed as in (31) and a new restart is used. The process is repeated until the norm of the restarted residual becomes smaller than $\mathrm{tol}_{\mathrm{cvg}}$. We summarize this discussion in the following algorithm, which will be referred to as Low-Rank Krylov Squared Smith (LRKSS).

A few comments are in order. Algorithm LRKSS computes $Q_1^E$ and $Q_1^F$ from $QR$ decompositions of $E$ and $F$ and applies a variant Algorithm 1 to $A$ starting with $Q_1^E$ and to $B$ starting with $Q_1^F$. In fact, our block Arnoldi implementation is based on Ruhe's version with elimination, see [17, p.197] or [16, Algorithm 6.1]. Then, at each iteration $j$, the matrices $Q_{j+1}^E$, $Q_{j+1}^F$, $\underline{\mathbb{H}}_j^E$ and $\underline{\mathbb{H}}_j^F$ are computed. If $j = 1$, then $Z_0^E = E$ and $Z_0^F = F$ and the residual norm $\|\Gamma_0\|$ is computed as in Proposition 5. Else, if $j$ is a power of 2, then the reduced SVDs are computed as in (21)–(23) by eliminating the same number of singular values which are smaller than $\mathrm{tol}_{\mathrm{svd}}$. The factors $Z_k^E$ and $Z_k^F$ are then computed as well as the corresponding residual norm. If the residual norm is larger than $\mathrm{tol}_{\mathrm{cvg}}$ and the size of the Krylov bases reaches its maximum $m_{\max}$, the factors $Z_k^E$ and $Z_k^F$ are updated and the algorithm is restarted with new matrices $Q_1^E$ and $Q_1^F$ obtained from a reduced SVD of the residual as in (32).

It seems difficult to find optimal choices of $\mathrm{tol}_{\mathrm{svd}}$ and $\mathrm{tol}_{\mathrm{cvg}}$. We notice, however, that $\mathrm{tol}_{\mathrm{cvg}}$ should not be chosen too small compared to $\mathrm{tol}_{\mathrm{svd}}$ for the SVDs in (21) and (22) would be such that $\|S_k^E\| < \mathrm{tol}_{\mathrm{svd}}$ and $\|S_k^F\| < \mathrm{tol}_{\mathrm{svd}}$, which leads to no improvement of the approximate solution.

With this way of restarting, the proposed squared Smith version can be applied to large matrices and this was our primary objective. However, the quadratic convergence of the original squared Smith method may be lost. The purpose of the next section is to accelerate the convergence by replacing equation (1) with an equivalent one with matrices having smaller spectral radii.

## 5 ADI iteration

The ADI method is an important iterative process for solving Lyapunov and Sylvester equations [24, 25, 5, 6, 15, 18, 4, 3, 13, 7]. An ADI iteration suited for equation (1) is proposed in [6], see also [4]. It is given, for $i = 0, 1, \ldots$, by

$$X_{i+\frac{1}{2}}(I - \delta_i B^T) = (A - \delta_i I)X_i B^T + EF^T,$$

$$(I - \eta_i A)X_{i+1} = AX_{i+\frac{1}{2}}(B^T - \eta_i I) + EF^T,$$

where $X_0$ is an initial approximate solution of (1) and $\mu_i$ and $\eta_i$ are parameters chosen to accelerate the convergence. Eliminating $X_{i+\frac{1}{2}}$ from these two equations and rearranging the terms, we obtain

$$
\begin{aligned}
X_{i+1} &= (I - \eta_i A)^{-1} A (A - \delta_i I) X_i B^T (I - \delta_i B^T)^{-1}(B^T - \eta_i I) + \\
&\quad (I - \eta_i A)^{-1}((1 - \delta_i \eta_i)AEF^T B^T)(I - \delta_i B^T)^{-1} + EF^T.
\end{aligned}
$$

14

**Algorithm 2** LRKSS

INPUT: $A, B \in \mathbb{R}^{n,n}, E, F \in \mathbb{R}^{n,p}$, an integer $m_{\mathrm{max}}$, tolerances $\mathrm{tol}_{\mathrm{svd}}$, $\mathrm{tol}_{\mathrm{cvg}}$, and an initial value of $\|\Gamma\| > \mathrm{tol}_{\mathrm{cvg}}$.

OUTPUT: Approximate solution to (1) in factored form $Z_s^E \left(Z_s^F\right)^T$

1: $QR$ factorize $E = Q_1^E R_1^E$ and $F = Q_1^F R_1^F$.
2: Set $U^E = I,\ U^F = I,\ V^E = I,\ V^F = I,\ S_1^E = R_1^E,\ S_1^F = R_1^F,\ Z_s^E = [\ ],\ Z_s^F = [\ ],$
$\quad p^E = \dim(Q_1^E),\ p^F = \dim(Q_1^F),$
3: $j = 0,\ iter = 0,\ rst = 0.$
4: **while** $(\|\Gamma\| > \mathrm{tol}_{\mathrm{cvg}})$ **do**
5: $\quad j := j + 1$
6: $\quad$ Update $\underline{\mathbb{H}}_j^E,\ \underline{\mathbb{H}}_j^F,\ \mathbb{Q}_{j+1}^E$ and $\mathbb{Q}_{j+1}^F$.
7: $\quad$ **if** $j = 2^k$ **then**
8: $\quad\quad$ **if** $k = 0$ **then**
9: $\quad\quad\quad Z^E = E,\ Z^F = F,\ \|\Gamma\| = \|\left(\mathbb{H}_1^E R_1^E\right)\left(\mathbb{H}_1^F R_1^F\right)^T\|$
10: $\quad\quad$ **else**
11: $\quad\quad\quad$ compute the reduced SVDs

$$U^E S^E (V^E)^T := \left( \begin{pmatrix} U^E S^E (V^E)^T V^F \\ 0 \end{pmatrix}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^E U^E S^E (V^E)^T V^F \right)$$

$\quad\quad\quad$ and

$$U^F S^F (V^F)^T := \left( \begin{pmatrix} U^F S^F \\ 0 \end{pmatrix}, \prod_{j=2^k-1}^{2^{k-1}} \mathbb{H}_j^F U^F S^F \right)$$

$\quad\quad\quad$ by eliminating the same number of singular values of $S^E$ and $S^F$ which are less than $\mathrm{tol}_{\mathrm{svd}}$.
12: $\quad\quad\quad Z^E = \mathbb{Q}_j^E U^E S^E (V^E)^T V^F, \quad Z^F = \mathbb{Q}_j^F U^F S^F$
13: $\quad\quad\quad \|\Gamma\| = \left\| \begin{pmatrix} R_1^E (R_1^F)^T & 0 \\ 0 & 0 \end{pmatrix} + \left(\underline{\mathbb{H}}_j^E U^E S^E (V^E)^T V^F\right)\left(\underline{\mathbb{H}}_j^F U^F S^F\right)^T \right.$
$\quad\quad\quad\quad \left. - \begin{pmatrix} U^E S^E (V^E)^T V^F (U^F)^T (S^F)^T & 0 \\ 0 & 0 \end{pmatrix} \right\|$
14: $\quad\quad$ **end if**
15: $\quad\quad iter := iter + 1$
16: $\quad$ **end if**
17: $\quad$ **if** $(j+1)p^E > m_{\mathrm{max}}$ or $\|\Gamma\| \le \mathrm{tol}_{\mathrm{cvg}}$ **then**
18: $\quad\quad$ Set $Z_s^E := (Z_s^E, Z^E),\ Z_s^F := (Z_s^F, Z^F).$
19: $\quad\quad$ Compute the reduced SVD

$$U^E S^E (V^E)^T \ := \ \begin{pmatrix} R_1^E (R_1^F)^T & 0 \\ 0 & 0 \end{pmatrix} + \left(\underline{\mathbb{H}}_j^E U^E S^E (V^E)^T V^F\right)\left(\underline{\mathbb{H}}_j^F U^F S^F\right)^T -$$
$$\begin{pmatrix} U^E S^E (V^E)^T V^F (U^F)^T (S^F)^T & 0 \\ 0 & 0 \end{pmatrix}$$

$\quad\quad$ with $\sigma_{\mathrm{min}}(S) > tol_{svd}.$
20: $\quad\quad$ Set $R_1^E = S^{\frac{1}{2}},\ R_1^F = S^{\frac{1}{2}}, Q_1^E = \mathbb{Q}_{j+1}^E U,\ Q_1^F = \mathbb{Q}_{j+1}^F V.$
21: $\quad\quad$ Set $U^E = I,\ U^F = I,\ V^E = I,\ V^F = I,\ S_1^E = R_1^E,\ S_1^F = R_1^F,\ j = 0,$
$\quad\quad rst := rst + 1,\ p^E = \dim(Q_1^E),\ p^F = \dim(Q_1^F).$
22: $\quad$ **end if**
23: **end while**

Let
$$\mathcal{A}_i = (I - \eta_i A)^{-1} A \ (A - \delta_i I), \quad \mathcal{B}_i = (I - \delta_i B)^{-1} B \ (B - \eta_i I),$$

$$\mathcal{E}_i = (E \ , \ (I - \eta_i A)^{-1} AE \ \sqrt{1 - \delta_i \eta_i}), \quad \mathcal{F}_i = (F \ , \ (I - \delta_i B)^{-1} BF \ \sqrt{1 - \delta_i \eta_i}).$$

Then the sequence $\left(X_i\right)_{i \geq 0}$ satisfies the iteration

$$X_{i+1} = \mathcal{A}_i X_i \mathcal{B}_i^T + \mathcal{E}_i \mathcal{F}_i^T. \tag{35}$$

A straightforward calculation shows that the solution $X$ is a fixed point of this iteration and hence the error is given by

$$X_{i+1} - X = \mathcal{A}_i \left(X_i - X\right) \mathcal{B}_i^T.$$

A repetition of this iteration gives

$$X_{i+1} - X = \left(\Pi_{j=0}^i \mathcal{A}_i\right) \left(X_0 - X\right) \left(\Pi_{j=0}^i \mathcal{B}_i\right)^T.$$

The convergence $X_{i+1} - X \to 0$ is fast if the spectral radii of $\Pi_{j=0}^i \mathcal{A}_j$ and $\Pi_{j=0}^i \mathcal{B}_j$ are as small as possible. Ideally, this will be the case if the parameters $\mu_i$ and $\eta_i$, $i = 0, 1, \ldots, i$, are chosen to satisfy

$$\begin{array}{cc} \min & \max \\ \delta_i \in \mathbb{C} & \lambda \in \Lambda(A) \\ \eta_i \in \mathbb{C} & \mu \in \Lambda(B) \end{array} \left| \Pi_{j=0}^i \frac{\lambda(\lambda - \delta_j)\mu(\mu - \eta_j)}{(1 - \eta_j \lambda)(1 - \delta_j \mu)} \right|. \tag{36}$$

However, this problem is hard to solve and computationally expensive. Since we are only interested in parameters that help reduce the spectral radii, we will consider only two parameters $\mu$ and $\eta$ which approximate the minimax problem

$$\begin{array}{cc} \min & \max \\ \delta \in \mathbb{C} & \lambda \in \Lambda(\tilde{A}) \\ \eta \in \mathbb{C} & \mu \in \Lambda(\tilde{B})| \end{array} \frac{|\lambda(\lambda - \delta)\mu(\mu - \eta)|}{|(1 - \eta\lambda)(1 - \delta\mu)|}, \tag{37}$$

where the sets $\Lambda(\tilde{A})$ and $\Lambda(\tilde{B})$ contain a few smallest and/or largest approximate eigenvalues of $A$ and $B$. The numerical tests show that this simple choice gives a satisfactory acceleration.

Once $\delta$ and $\eta$ are computed, the matrices $A$, $B$, $E$ and $F$ of equation (1) are replaced by

$$\mathcal{A} = (I - \eta A)^{-1} A(A - \delta I), \quad \mathcal{B} = (I - \delta B)^{-1} B(B - \eta I) \tag{38}$$

$$\mathcal{E} = \left(E, (I - \eta A)^{-1} AE\sqrt{1 - \delta\eta}\right), \quad \mathcal{F} = \left(F, (I - \delta B)^{-1} BF\sqrt{1 - \delta\eta}\right) \tag{39}$$

on which algorithm LRKSS is applied.

# 6 Numerical Tests

We present numerical tests to illustrate the key points crucial for the convergence of algorithm LRKSS and its ADI acceleration.

**Test 1:** This test shows the convergence behavior when the spectral radii of $A$ and $B$ approach 1. The matrices $A$ and $B$ are $n \times n$, Toeplitz tridiagonal. $A$ has $-\alpha$, 0 and $+\alpha$ respectively on its subdiagonal, diagonal and superdiagonal and $B$ has the same structure with $-\beta$, 0 and $+\beta$, where $\alpha$ and $\beta$ are positive parameters to be varied. The matrix $E$ is $n \times 2$ formed by the first two vectors of the canonical basis, that is, $E_{11} = E_{22} = 1$ and zero elsewhere, and $F = -E$. The eigenvalues of $A$ and $B$ are given by $\lambda_j = 2i\alpha \cos \frac{\pi j}{n+1}$ and $\mu_j = 2i\beta \cos \frac{\pi j}{n+1}$, $j = 1, ..., n$. For large $n$ we see that $\rho(A) \approx 2\alpha$ and $\rho(B) \approx 2\beta$.

Table 1 shows the results obtained with $n = 10^3$, $\text{tol}_\text{cvg} = \text{tol}_\text{svd} = 10^{-10}$ and different values of $\alpha$, $\beta$ and the maximum dimension of the Krylov spaces $m_\text{max}$. The table also indicates the corresponding spectral radii $\rho(A)$ and $\rho(B)$, residual norms, total number of iterations and restarts. As expected, the closer the spectral radii get to 1, the slower the convergence is. Also, note that the restart, while it remedies the problem of storage requirements and computational cost, slows down the convergence. The numbers of restarts and iterations are almost doubled when the dimension of the Krylov spaces is divided by 2.

| $\alpha$ $\beta$ | $\rho(A)$ $\rho(B)$ | $m_\text{max}$ | res.norm | iter | rst |
|---|---|---|---|---|---|
| 0.45 0.445 | 0.9 0.89 | 32 | $5.96 \times 10^{-11}$ | 20 | 4 |
| | | 64 | $5.96 \times 10^{-11}$ | 14 | 2 |
| | | 128 | $5.96 \times 10^{-11}$ | 10 | 1 |
| 0.499 0.495 | 0.998 0.99 | 32 | $9.94 \times 10^{-11}$ | 268 | 66 |
| | | 64 | $8.24 \times 10^{-11}$ | 171 | 33 |
| | | 128 | $5.35 \times 10^{-11}$ | 102 | 16 |
| 0.4999 0.499 | 0.9998 0.998 | 32 | $9.93 \times 10^{-11}$ | 1205 | 296 |
| | | 64 | $9.93 \times 10^{-11}$ | 753 | 148 |
| | | 128 | $9.93 \times 10^{-11}$ | 452 | 74 |

Table 1: Results of LRKSS with different values of $\alpha$, $\beta$ and $m_\text{max}$ (Test 1)

**Test 2:** We use the matrices as in the previous test with fixed $\alpha = 0.499$, $\beta = 0.495$, $m_\text{max} = 64$ and consider three values of $n$, $n = 10^3$, $n = 10^4$, and $n = 10^5$. The corresponding spectral radii are almost the same $\left(\rho(A) \approx 2\alpha, \rho(B) \approx 2\beta\right)$ meaning that the convergence behavior is almost the same, see Proposition 3. For the three tests, the numbers of iterations and restarts are respectively 171 and 33. The convergence behaviors are shown in Figure 1.
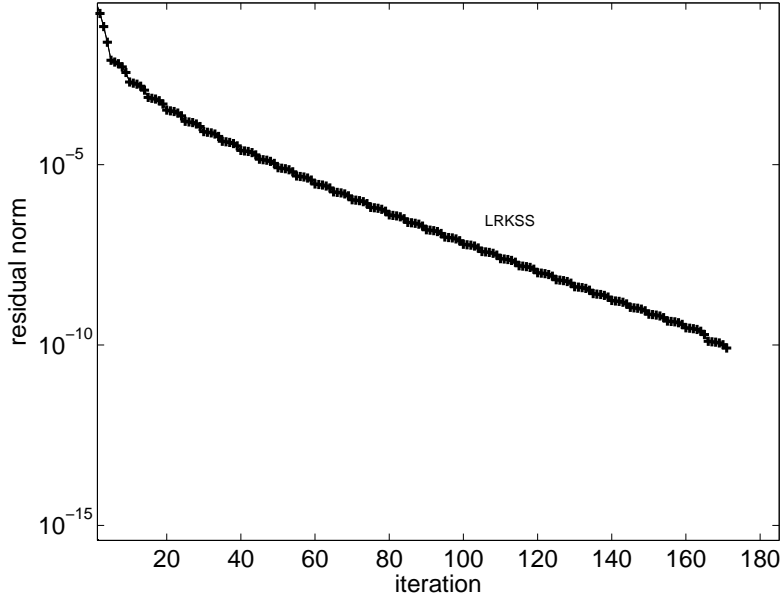
Figure 1: Residual norms vs the number of iterations (Test 2, $n = 10^p$, $p = 3, 4, 5$)

**Test 3:** We consider now Test 1 with ADI acceleration. The parameters $\delta$ and $\eta$ of the minimax problem (37) are obtained from the 10 largest (in modulus) eigenvalues of $A$ and $B$. The parameters $\delta$ and $\eta$ are computed by the MATLAB function `fminsearch` as proposed in [4]. The linear systems in (38) and (39) are solved by GMRES with restart value 20 and tolerance $10^{-10}$. The results are summarized in Table 2. An improvement can be noticed compared to Table 1. Note that the parameters $\eta$ and $\mu$ are close to zero, which means that the matrices $\mathcal{A}$, $\mathcal{B}$, $\mathcal{E}$ and $\mathcal{F}$ in (38) and (39) are close to $A^2$, $B^2$, $(E, AE)$ and $(F, BF)$.

| $\alpha$ $\beta$ | $\delta$ $\eta$ | $\rho(\mathcal{A})$ $\rho(\mathcal{B})$ | $m_{\max}$ | res.norm | iter | rst |
|---|---|---|---|---|---|---|
| 0.45 | $-7.877\ 10^{-9}$ | 0.8099 | 32 | $7.15 \times 10^{-11}$ | 13 | 3 |
| 0.445 | $-7.877\ 10^{-9}$ | 0.7921 | 64 | $7.15 \times 10^{-11}$ | 11 | 2 |
| | | | 128 | $7.15 \times 10^{-11}$ | 8 | 1 |
| 0.499 | $-4.3315\ 10^{-8}$ | 0.9959 | 32 | $8.30 \times 10^{-11}$ | 159 | 43 |
| 0.495 | $3.5532\ 10^{-8}$ | 0.9801 | 64 | $9.74 \times 10^{-11}$ | 102 | 22 |
| | | | 128 | $6.70 \times 10^{-11}$ | 66 | 11 |
| 0.4999 | $1.5357\ 10^{-7}$ | 0.9995 | 32 | $9.76 \times 10^{-11}$ | 670 | 173 |
| 0.499 | $-1.1153\ 10^{-7}$ | 0.99959 | 64 | $9.76 \times 10^{-11}$ | 424 | 86 |
| | | | 128 | $9.76 \times 10^{-11}$ | 256 | 42 |

Table 2: Results of LRKSS and ADI with different values of $\alpha$, $\beta$ and $m_{\max}$ (Test 3)

18

**Test 4:**  We consider now the equivalent equation

$$X - A^2 X (B^2)^T = (E, AE)(F, AF)^T$$

on which we apply ADI iterations and algorithm LRKSS, $n = 10^3$, $\text{tol}_{\text{svd}} = \text{tol}_{\text{cvg}} = 10^{-10}$. GMRES is used with the same parameters as in Test 3. Note that the parameters $\eta$ and $\mu$ allow now a significant improvement compared to the previous results, see Table 3.

| $\alpha$ $\beta$ | $\delta$ $\eta$ | $\rho(\mathcal{A})$ $\rho(\mathcal{B})$ | $m_{\max}$ | res.norm | iter | rst |
|---|---|---|---|---|---|---|
| 0.45 | $-8.0986\ 10^{-1}$ | 0.2565 | 32 | $9.43 \times 10^{-11}$ | 4 | 0 |
| 0.445 | $-7.9195\ 10^{-1}$ | 0.2453 | 64 | $9.43 \times 10^{-11}$ | 3 | 0 |
| | | | 128 | $9.43 \times 10^{-11}$ | 3 | 0 |
| 0.499 | $-9.9582\ 10^{-1}$ | 0.7427 | 32 | $2.10 \times 10^{-11}$ | 16 | 6 |
| 0.495 | $-9.7992\ 10^{-1}$ | 0.7192 | 64 | $1.62 \times 10^{-12}$ | 13 | 3 |
| | | | 128 | $7.59 \times 10^{-11}$ | 9 | 2 |
| 0.4999 | $-9.9583\ 10^{-1}$ | 0.8742 | 32 | $4.04 \times 10^{-11}$ | 31 | 12 |
| 0.499 | $-9.9942\ 10^{-1}$ | 0.8679 | 64 | $7.12 \times 10^{-11}$ | 24 | 7 |
| | | | 128 | $1.30 \times 10^{-11}$ | 17 | 3 |

Table 3: Results of LRKSS and ADI with different values of $\alpha$, $\beta$ and $m_{\max}$ (Test 4)

Figure 2 draws the convergence behaviors when $n = 10^3$, $n = 10^4$ and $n = 10^5$. For the three cases, $\alpha = 0.499$, $\beta = 0.495$ and $m_{\max} = 64$. The three indistinguishable curves in this figure show that the convergence behavior is the same.
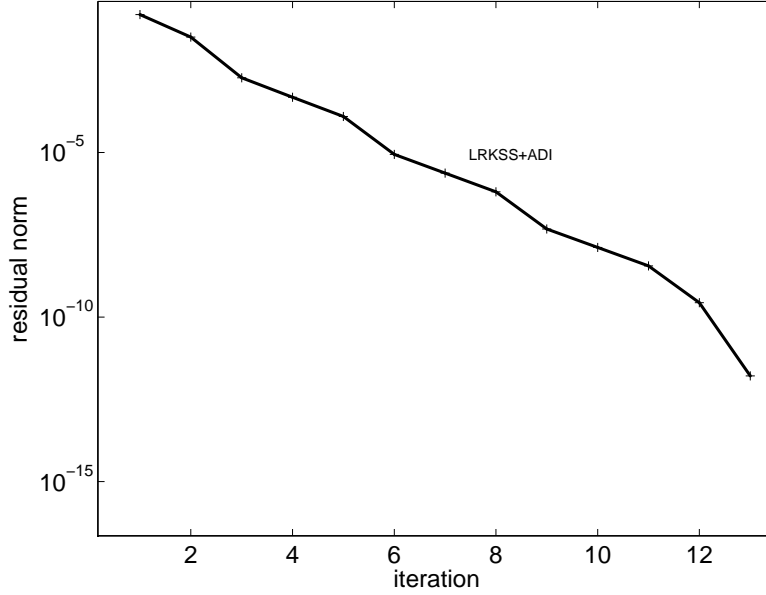
Figure 2: Residual norms vs the number of iterations, (Test 4, $n = 10^p$, $p = 3, 4, 5$)

**Test 5:** In this test we solve the equation

$$X - AXA = EF^T,$$

where $A = Q^T \hat{A} Q$, $Q$ is an orthogonal matrix constructed with the MATLAB function `orth`, $\hat{A} = \text{diag}\left(\left(0.999 \, e^{i\frac{\pi}{n}}\right)^k, \; 1 \leq k \leq n\right)$, and the matrices $E$ and $F$ are the same as in the previous tests.

Taking $\eta = \delta$ in (37) and omitting the spectral part due to $B$, the minimax problem simplifies to

$$\min_{|\delta|<1} \max_{\lambda \in \Lambda(\tilde{A})} \frac{|\lambda(\lambda - \delta)|}{|(1 - \delta\lambda)|}. \tag{40}$$

The matrices in (38) and (39) become

$$\mathcal{A} = (I - \delta A)^{-1} A (A - \delta I), \tag{41}$$

$$\mathcal{E} = \left(E, (I - \delta A)^{-1} AE \sqrt{1 - \delta^2}\right), \quad \mathcal{F} = \left(F, (I - \delta A^T)^{-1} A^T F \sqrt{1 - \delta^2}\right) \tag{42}$$

Taking $n = 10^3$ we obtain $\delta = 9.8280 \times 10^{-1}$, $\rho(\mathcal{A}) = 8.9335 \times 10^{-1}$. The parameters $\text{tol}_{\text{svd}}$ and $\text{tol}_{\text{cvg}}$ are fixed at $10^{-10}$ and $m_{\max} = 32$. Figure 3 on the left shows the convergence of LRKSS with and without ADI preconditioning. The figure on the right shows the singular values of the exact and computed solutions. The smallest singular

values differ by a factor of order $10^{-10} (= \text{tol}_{\text{svd}})$. Figure 4 shows the convergence with different values of $m_{\max}$ and confirms again that the larger the values of $m_{\max}$ is, the smaller the number of iterations.
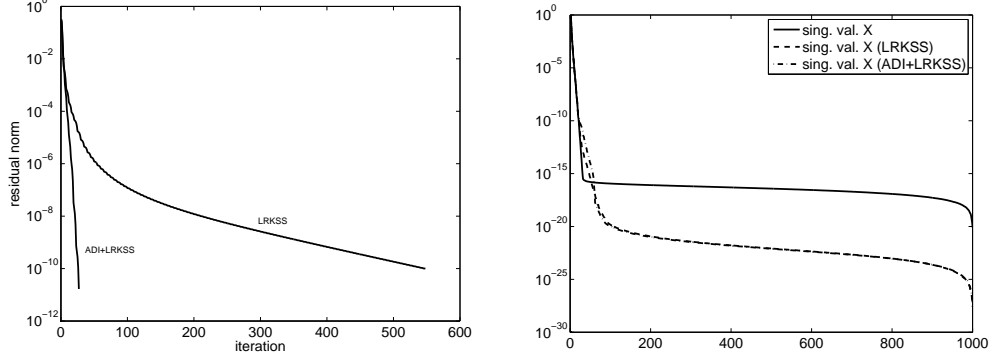


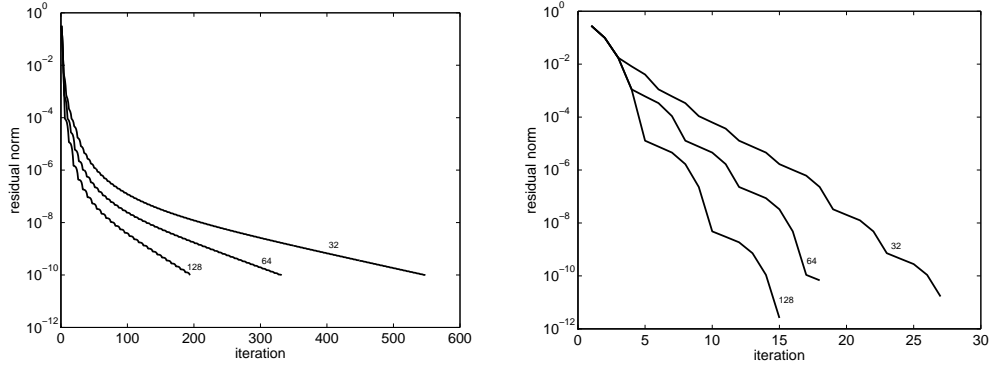Figure 3: Residual norms and singular values of the exact and the computed solutions (Test 5)



Figure 4: Convergence behaviors of LKRSS (left) and LKRSS+ADI(right) with different values of $m_{\max}$ (Test 5)

## 7 Conclusions

The main purpose of this work was to show one way of adapting the Squared Smith method to large-scale Stein equations. The adaptation requires the use of Krylov spaces to build approximations of the squared Smith iterates in low-rank factors. As expected, the quadratic convergence in the original squared Smith algorithm is not maintained, but the association with the proposed adaptation with a simple version of the ADI iteration as a preconditioner allows a great acceleration of the convergence.

21

This is consistent with the numerical results in [13, 7] where the "optimal" number of ADI iterations is less than 3 in [13] and around 4 in [7]. The acceleration depends largely on the ADI parameters and to a lesser extend on the other parameters of the algorithm. Improvements can still be made if these parameters can be chosen in a cheap and nearly optimal way.

# References

[1] Björck A. Numerical Methods for Least Squares Problems. SIAM: Philadelphia, 1996.

[2] Antoulas AC. Approximation of Large-Scale Dynamical Systems, Advances in Design and Control. SIAM: Philadelphia, 2005.

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

[4] Benner P, Faßbender H. On the numerical solution of large-scale sparse discrete-time Riccati equations. Advances in Computational Mathematics. 2011; 35: 119--147.

[5] Calvetti D, Reichel L. Application of ADI iterative methods to the restoration of noisy images. SIAM J. Matrix Anal. Appl. 1996; 17: 165--186.

[6] Calvetti D, Levenberg N, Reichel L. Iterative methods for $X - AXB = C$. J. Comput. Appl. Math. 1997; 86: 73--101.

[7] Damm T. Direct methods and ADI-preconditioned Krylov subspace methods for generalized Lyapunov equations. Numer. Linear Algebra Appl. 2008; 15: 853--871.

[8] El Guennouni A, Jbilou K, Riquet J. Block Krylov subspace methods for solving large Sylvester equations. Numerical Algorithms. 2001; 29: 75--96.

[9] Gajič Z, Qureshi MTJ. Lyapunov matrix equation in system stability and control. Mineola, N. Y., Dover, 2008.

[10] Horn RA, Johnson CR. Topics in Matrix Analysis. Cambridge Universtity Press, Cambridge, 1990.

[11] Hu DY, Reichel L. Krylov subspace methods for the Sylvester equations. Linear Algebra Appl. 1992; 174: 283--314.

[12] Ionescu V, Oara C, Weiss M. Generalized Riccati theroy and robust control. A Popov function approach, John Wiley & Sons, Ltd., Chichester, 1999.

[13] Jbilou K. ADI preconditioned Krylov methods for large Lyapunov matrix equations. Linear Algebra Appl., 2010; 432: 2473--2485.

[14] Lancaster P, Rodman L. Algebraic Riccati Equations. Clarendon Press, Oxford, 1995.

[15] Penzl T. A cyclic low-rank Smith method for large sparse Lyapunov equations. SIAM J. Sci. Comput.2000; 21: 1401--1418.

[16] Robbé M, Sadkane M. A convergence analysis of GMRES and FOM methods for Sylvester equations. Numerical Algorithms. 2002; 30: 71--89.

[17] Saad Y. Iterative methods for sparse linear systems, 2nd ed. SIAM: Philadelphia, 2003.

[18] Sabino J. Solution of large-scale Lyapunov equations via the block modified Smith method. PhD thesis, Rice University, 2006.

[19] Sadkane M. A low-rank squared Smith method for large-scale discrete-time Lyapunov equations. Linear Algebra Appl. 2012; 436: 2807--2827.

[20] Simoncini V. On the numerical Solution of $AX - XB = C$. BIT. 1996; 36: 814--830.

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

[22] Stewart GW, Sun JG. Matrix Perturbation Theory. Academic Press, San Diego, CA, 1990.

[23] Trefethen LN, Embree M. Spectra and pseudospectra. The behavior of nonnormal matrices and operators. Princeton University Press, Princeton, 2005.

[24] Wachspress EL. Extended application of alternating direction implicit iteration model problem theory. Journal of SIAM. 1963; 11: 994--1016.

[25] Wachspress EL. Iterative solution of the Lyapunov matrix equations. Appl. Math. Letters. 1988; 107: 87--90.