

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/222552646>

# A structure-preserving doubling algorithm for continuous-time algebraic Riccati equations

Article in *Linear Algebra and its Applications* · February 2005

DOI: 10.1016/j.laa.2004.10.010

CITATIONS

106

READS

144

3 authors, including:



[Eric Chu](#)

Monash University (Australia)

131 PUBLICATIONS 2,056 CITATIONS

[SEE PROFILE](#)



[Wen-Wei Lin](#)

National Chiao Tung University

224 PUBLICATIONS 3,490 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Numerical methods for large matrix eigenproblems [View project](#)



3D Maxwell Eigenvalue Problems & Photonic Crystal [View project](#)



ELSEVIER

Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

SCIENCE @ DIRECT®

LINEAR ALGEBRA  
AND ITS  
APPLICATIONS

Linear Algebra and its Applications 396 (2005) 55–80

[www.elsevier.com/locate/laa](http://www.elsevier.com/locate/laa)

## A structure-preserving doubling algorithm for continuous-time algebraic Riccati equations

E.K.-W. Chu<sup>a,\*</sup>, H.-Y. Fan<sup>b</sup>, W.-W. Lin<sup>b</sup>

<sup>a</sup>*School of Mathematical Sciences, Monash University, Building 28, VIC 3800, Australia*

<sup>b</sup>*Department of Finance, Yuanpei Institute of Science and Technology, Hsinchu 300, Taiwan*

Received 12 June 2003; accepted 6 October 2004

Submitted by R. Byers

---

### Abstract

Continuous-time algebraic Riccati equations (CAREs) can be transformed, *à la* Cayley, to discrete-time algebraic Riccati equations (DAREs). The efficient structure-preserving doubling algorithm (SDA) for DAREs, from [E.K.-W. Chu, H.-Y. Fan, W.-W. Lin, A structure-preserving doubling algorithm for periodic discrete-time algebraic Riccati equations, preprint 2002-28, NCTS, National Tsing Hua University, Hsinchu 300, Taiwan, 2003; E.K.-W. Chu, H.-Y. Fan, W.-W. Lin, C.-S. Wang, A structure-preserving doubling algorithm for periodic discrete-time algebraic Riccati equations, preprint 2002-18, NCTS, National Tsing Hua University, Hsinchu 300, Taiwan, 2003], can then be applied. In this paper, we develop the structure-preserving doubling algorithm from a new point of view and show its quadratic convergence under assumptions which are weaker than stabilizability and detectability, as well as practical issues involved in the application of the SDA to CAREs. A modified version of the SDA, developed for DAREs with a “doubly symmetric” structure, is also presented. Extensive numerical results show that our approach is efficient and competitive.

© 2004 Elsevier Inc. All rights reserved.

*AMS classification:* 93B50; 93B52; 93C05; 93D15

*Keywords:* Cayley transform; Continuous-time algebraic Riccati equation; Doubling algorithm; Matrix sign function; Structure-preserving

---

\* Corresponding author.

*E-mail addresses:* [eric.chu@sci.monash.edu.au](mailto:eric.chu@sci.monash.edu.au) (E.K.-W. Chu), [hyfan@mail.yust.edu.tw](mailto:hyfan@mail.yust.edu.tw) (H.-Y. Fan), [wwlin@am.nthu.edu.tw](mailto:wwlin@am.nthu.edu.tw) (W.-W. Lin).

## 1. Introduction

In this paper we investigate a structure-preserving doubling algorithm [24,37] for the computation of the symmetric positive semi-definite (s.p.s.d.) solution  $X$  (i.e.  $X \geq 0$ ) to the continuous-time algebraic Riccati equation (CARE):

$$-XGX + A^T X + XA + H = 0, \quad (1)$$

where  $A \in \mathbb{R}^{n \times n}$ ,  $X \in \mathbb{R}^{n \times n}$ ,  $R \in \mathbb{R}^{m \times m}$  is symmetric positive definite (or s.p.d.; i.e.  $R > 0$ ),  $G = BR^{-1}B^T \geq 0$  and  $H = C^T C \geq 0$  with  $B \in \mathbb{R}^{n \times m}$  and  $C^T \in \mathbb{R}^{n \times p}$  being of full column rank.

Eq. (1) arises frequently in solving the continuous-time linear optimal control problem:

$$\min_u J = \frac{1}{2} \int_0^\infty (x^T C^T C x + u^T R u) dt \quad \text{subject to } \dot{x} = Ax + Bu. \quad (2)$$

The optimal feedback control  $u^*$  for (2) is given by

$$u^* = -R^{-1}B^T X x, \quad (4)$$

where  $X$  is the s.p.s.d. solution to the CARE (1). We assume that the pair  $(A, B)$  is stabilizable (S) (i.e. if  $w^T B = 0$  and  $w^T A = \lambda w^T$  for some  $\lambda \in \mathbb{C}$ , then  $\text{Re}(\lambda) < 0$  or  $w = 0$ ) and that the pair  $(A, C)$  is detectable (D) (i.e.  $(A^T, C^T)$  is stabilizable). Under assumptions (S) and (D), the CARE (1) has been proved to possess a unique s.p.s.d. solution [39].

Consider the  $2n \times 2n$  Hamiltonian matrix  $\mathcal{H}$  associated with the CARE (1):

$$\mathcal{H} = \begin{bmatrix} A & -G \\ -H & -A^T \end{bmatrix}, \quad (5)$$

which satisfies

$$\mathcal{H}J = -J\mathcal{H}^T, \quad J = \begin{bmatrix} 0 & I_n \\ -I_n & 0 \end{bmatrix}$$

with  $I_n$  denoting the identity matrix of order  $n$ . By (5), the CARE (1) can be written as

$$\mathcal{H} \begin{bmatrix} I \\ X \end{bmatrix} = \begin{bmatrix} I \\ X \end{bmatrix} \Phi, \quad (6)$$

where  $\Phi \in \mathbb{R}^{n \times n}$  and the spectrum  $\sigma(\Phi)$  is on the stable left half plane  $\mathbb{C}_-$ . Under assumptions (S) and (D), the Hamiltonian matrix  $\mathcal{H}$  has exactly  $n$  eigenvalues on  $\mathbb{C}_-$ . If the columns of  $[X_1^T, X_2^T]^T$  span the stable invariant subspace of  $\mathcal{H}$ , then  $X_1$  is nonsingular and  $X = X_2 X_1^{-1} \geq 0$  solves the CARE (1) (see, e.g., [39,44]).

A numerically backward stable algorithm `care`, proposed by Laub [39], computes  $X$  by applying the QR algorithm with reordering [4,16,48] to the eigenvalue problem  $\mathcal{H}x = \lambda x$ . Unfortunately, the QR algorithm preserves neither the Hamiltonian

structure of  $\mathcal{H}$  nor the associated splitting of eigenvalues. A structure-preserving algorithm has been proposed by Ammar and Mehrmann [1] which utilizes orthogonal symplectic transformations in computing a basis for the stable invariant subspace of  $\mathcal{H}$ . A stable symplectic orthogonal method has been suggested by Byers [19] but applied only to systems with single input or output. Many iterative methods have been suggested for solving CAREs over the past 20 years. Newton's method has been applied in extensive literature [28,31,38,41,47]. A defect correction method for modifying an approximate solution has also been proposed by Mehrmann and Tan [43]. These methods require a good starting approximate solution, and can therefore be regarded as iterative refinement methods, to be combined with other direct methods (see Bunse-Gerstner et al. [17,18] or Mehrmann [41] for details). The structure-preserving matrix sign function methods (MSGM) [7,11–14,20,21,27,33,46] have been extended by Barraud [8,9] and Gardiner and Laub [29].

A class of methods, referred to as the doubling algorithms (DA), has attracted much interests in the 70s and 80s (see [2] and the references therein). These methods originate from the fixed-point iteration derived from the discrete-time algebraic Riccati equation (DARE):

$$X_{k+1} = \hat{A}^T X_k (I + \hat{G} X_k)^{-1} \hat{A} + \hat{H}.$$

Instead of producing the sequence  $\{X_k\}$ , doubling algorithms produce  $\{X_{2^k}\}$ . CAREs can be tackled after being transformed to DAREs via the Cayley transform. However, the convergence of the algorithm was proven only when  $\hat{A}$  is nonsingular [2], and for  $(\hat{A}, \hat{G}, \hat{H})$  which is stabilizable and detectable [36]. DAs were largely forgotten in the past decade. Recently, DAs have been revived for (periodic) DAREs, because of a better theoretical understanding. Stronger convergence results have been proved for  $(\hat{A}, \hat{G}, \hat{H})$  under weaker assumptions than stabilizability and detectability [24]. Superior numerical results, in comparison to state-of-the-art methods on a wide range of test problems, have been obtained because of the stronger structure-preserving properties and the superior operations count.

In this paper, we propose a doubling algorithm for CAREs. The CAREs are transformed to DAREs, with the corresponding Hamiltonian matrix transformed into a symplectic matrix pair by the Cayley transform. Nice convergence properties are inherited from the structure-preserving doubling algorithm (SDA) [24] applied to the corresponding DARE. The SDA preserves matrix pairs in SSF which is a stronger property than symplecticity. In the CARE setting, the matrix sign function methods preserve the Hamiltonian structure in  $\mathcal{H}$  while the SDA preserves, in each iterative step, the associated symplectic matrix pair  $(\hat{\mathcal{N}}, \hat{\mathcal{L}})$  in SSF. Although under the influence of numerical errors, the matrix pairs through the SDA retain their stabilizability, detectability as well as eigenstructures (with exactly half of the spectrum being stable; see details in [24]). This stronger structure-preserving property is its main strength and the reason of its accuracy. In Section 4, a modified version of the SDA (SDA\_m) is developed, for “doubly symmetric” DAREs, where  $\hat{A}, \hat{G} = \hat{H}$  are symmetric and persymmetric. The SDA\_m preserves the symplectic and doubly

symmetric structures of the DARE, resulting in better accuracy than the SDA. We have extensively tested the SDA against the MSGM and care. Numerical results showed that the doubling algorithm for CAREs is competitive and promising.

Finally, it is important to stress that matrix sign functions can be applied to more general Hamiltonian matrices in other applications, such as those from  $H_\infty$  control with  $G$  and  $H$  being indefinite. A scaling strategy [21] may also accelerate its convergence. Also, the SDA requires the transformation of the CARE by the Cayley transform, which requires the estimation of the parameter  $\gamma$  (see Section 3).

## 2. SDA and matrix sign function method

In this section we propose a structure-preserving doubling algorithm (SDA) for solving the CARE (1) based on the doubling algorithms in [24,37]. In addition, the well-known structure-preserving matrix sign function methods [7,11–14,20,21,27,33,46] are also reviewed from the point of view of preserving Hamiltonian structure.

Let  $\mathbf{H}$  be the set of  $2n \times 2n$  Hamiltonian matrices, i.e.,

$$\mathbf{H} = \left\{ \mathcal{H} \left| \mathcal{H} = \begin{bmatrix} A & -G \\ -H & -A^T \end{bmatrix}; A, H, G \in \mathbb{R}^{n \times n}; H, G \geq 0 \right. \right\}. \quad (7)$$

Note that if  $\mathcal{H} \in \mathbf{H}$  then  $\mathcal{H}J = -J\mathcal{H}^T$ . We call a  $2n \times 2n$  matrix pair  $(\mathcal{N}, \mathcal{L})$  symplectic if  $\mathcal{N}J\mathcal{N}^T = \mathcal{L}J\mathcal{L}^T$ . Let  $\mathbf{S}$  be the set of  $2n \times 2n$  symplectic matrix pairs in the standard symplectic form (SSF):

$$\mathbf{S} = \left\{ (\widehat{\mathcal{N}}, \widehat{\mathcal{L}}) \left| \widehat{\mathcal{L}} = \begin{bmatrix} I & \widehat{G} \\ 0 & \widehat{A}^T \end{bmatrix}, \widehat{\mathcal{N}} = \begin{bmatrix} \widehat{A} & 0 \\ -\widehat{H} & I \end{bmatrix}; \right. \right. \\ \left. \widehat{A}, \widehat{H}, \widehat{G} \in \mathbb{R}^{n \times n}; \widehat{G}, \widehat{H} \geq 0 \right\}. \quad (8)$$

It is easily seen that symplecticity is weaker than symplecticity in SSF. Our proposed algorithm preserves the stronger structure and gives rise to better numerical performance.

We shall show how the CARE (1), associated with the corresponding Hamiltonian matrix

$$\mathcal{H} \equiv \begin{bmatrix} A & -G \\ -H & -A^T \end{bmatrix} = \begin{bmatrix} A & -BR^{-1}B^T \\ -C^TC & -A^T \end{bmatrix} \in \mathbb{R}^{2n \times 2n}$$

can be transformed to an equivalent DARE.

By using the Cayley transform with some appropriate  $\gamma > 0$ , the Hamiltonian matrix  $\mathcal{H}$  can be transformed to a symplectic matrix pair  $(\mathcal{N}, \mathcal{L}) \equiv (\mathcal{H} + \gamma I, \mathcal{H} - \gamma I)$  [41,42]. In the following, we construct an equivalence transformation from  $(\mathcal{N}, \mathcal{L})$  to a symplectic matrix pair  $(\widehat{\mathcal{N}}, \widehat{\mathcal{L}}) \in \mathbf{S}$ .

Let

$$A_\gamma \equiv A - \gamma I, \quad \bar{A}_\gamma \equiv A + \gamma I.$$

Starting from

$$\mathcal{N} = \begin{bmatrix} \bar{A}_\gamma & -G \\ -H & -A_\gamma^T \end{bmatrix}, \quad \mathcal{L} = \begin{bmatrix} A_\gamma & -G \\ -H & -\bar{A}_\gamma^T \end{bmatrix},$$

we choose a  $\gamma > 0$  such that the matrices  $A_\gamma$  and  $A_\gamma + GA_\gamma^{-T}H$  are well-conditioned (see Section 3 later for details). To transform the symplectic matrix pair  $(\mathcal{N}, \mathcal{L})$  to  $(\widehat{\mathcal{N}}, \widehat{\mathcal{L}}) \in \mathbf{S}$ , let

$$T_1 \equiv \begin{bmatrix} A_\gamma^{-1} & 0 \\ HA_\gamma^{-1} & I \end{bmatrix}, \quad T_2 \equiv \begin{bmatrix} I & 0 \\ 0 & (-HA_\gamma^{-1}G - A_\gamma^T)^{-1} \end{bmatrix},$$

$$T_3 \equiv \begin{bmatrix} I & A_\gamma^{-1}G \\ 0 & I \end{bmatrix}.$$

Simple calculations produce

$$\begin{aligned} \widehat{\mathcal{N}} &= \begin{bmatrix} \widehat{A} & 0 \\ -\widehat{H} & I \end{bmatrix} = T_3 T_2 T_1 \mathcal{N} \\ &= T_3 T_2 \begin{bmatrix} A_\gamma^{-1} \bar{A}_\gamma & -A_\gamma^{-1}G \\ HA_\gamma^{-1} \bar{A}_\gamma - H & -HA_\gamma^{-1}G - A_\gamma^T \end{bmatrix} \\ &= T_3 \begin{bmatrix} A_\gamma^{-1} \bar{A}_\gamma & -A_\gamma^{-1}G \\ (-HA_\gamma^{-1}G - A_\gamma^T)^{-1} (HA_\gamma^{-1} \bar{A}_\gamma - H) & I \end{bmatrix} \end{aligned}$$

and

$$\begin{aligned} \widehat{\mathcal{L}} &= \begin{bmatrix} I & \widehat{G} \\ 0 & \widehat{A}^T \end{bmatrix} = T_3 T_2 T_1 \mathcal{L} \\ &= T_3 T_2 \begin{bmatrix} I & -A_\gamma^{-1}G \\ 0 & -HA_\gamma^{-1}G - \bar{A}_\gamma^T \end{bmatrix} \\ &= T_3 \begin{bmatrix} I & -A_\gamma^{-1}G \\ 0 & (-HA_\gamma^{-1}G - A_\gamma^T)^{-1} (-HA_\gamma^{-1}G - \bar{A}_\gamma^T) \end{bmatrix}, \end{aligned}$$

where

$$\begin{aligned} \widehat{A} &= (\bar{A}_\gamma + GA_\gamma^{-T}H)(A_\gamma + GA_\gamma^{-T}H)^{-1}, \\ \widehat{G} &= -A_\gamma^{-1}G + A_\gamma^{-1}G(A_\gamma^T + HA_\gamma^{-1}G)^{-1}(\bar{A}_\gamma^T + HA_\gamma^{-1}G), \\ \widehat{H} &= (A_\gamma^T + HA_\gamma^{-1}G)^{-1}(HA_\gamma^{-1}\bar{A}_\gamma - H). \end{aligned}$$

Note that  $\mathcal{L}^{-1}\mathcal{N} = \widehat{\mathcal{L}}^{-1}\widehat{\mathcal{N}}$ . Since  $\bar{A}_\gamma = A_\gamma + 2\gamma I$ , it follows that:

$$\widehat{A} = I + 2\gamma(A_\gamma + GA_\gamma^{-T}H)^{-1}, \quad (9)$$

$$\widehat{G} = 2\gamma A_\gamma^{-1} G (A_\gamma^T + H A_\gamma^{-1} G)^{-1} = 2\gamma A_\gamma^{-1} G A_\gamma^{-T} (I + H A_\gamma^{-1} G A_\gamma^{-T})^{-1}, \quad (10)$$

$$\widehat{H} = 2\gamma (A_\gamma^T + H A_\gamma^{-1} G)^{-1} H A_\gamma^{-1} = 2\gamma (I + A_\gamma^{-T} H A_\gamma^{-1} G)^{-1} A_\gamma^{-T} H A_\gamma^{-1}. \quad (11)$$

From (10), (11) and Lemma A.1 in Appendix A, we know that the matrices  $\widehat{G}$  and  $\widehat{H}$  are positive semi-definite. We thus obtain the desired symplectic matrix pair in SSF, i.e.,

$$(\widehat{\mathcal{N}}, \widehat{\mathcal{L}}) \equiv \left( \begin{bmatrix} \widehat{A} & 0 \\ -\widehat{H} & I \end{bmatrix}, \begin{bmatrix} I & \widehat{G} \\ 0 & \widehat{A}^T \end{bmatrix} \right) \in \mathbf{S},$$

where  $\widehat{A}$ ,  $\widehat{G}$  and  $\widehat{H}$  are given by (9)–(11). The DARE associated with the symplectic matrix pair  $(\widehat{\mathcal{N}}, \widehat{\mathcal{L}})$  in SSF is

$$X = \widehat{A}^T X (I + \widehat{G} X)^{-1} \widehat{A} + \widehat{H} \quad (12)$$

on which the efficient SDA [24] can be applied. Note that  $X$  is the unique s.p.s.d. solution to the above DARE as well as the CARE (1). Moreover, in Theorems 1 and 2 of [37], the pairs  $(\widehat{A}, \widehat{B})$  and  $(\widehat{A}, \widehat{C})$  are proven to be stabilizable and detectable, respectively, where the matrices  $\widehat{G} = \widehat{B} \widehat{B}^T$  and  $\widehat{H} = \widehat{C}^T \widehat{C}$  are full rank decompositions (FRD).

Using (9)–(11) to transform the CARE (1) to an equivalent DARE (12) with the associated symplectic matrix pair  $(\widehat{\mathcal{N}}, \widehat{\mathcal{L}})$  in SSF, the SDA in [24] can then be modified to the following algorithm for CAREs: (with  $\text{Im}$  denoting the imaginary axis).

### 2.1. Structure-preserving doubling algorithm (SDA):

Input:  $\mathcal{H} = \begin{bmatrix} A & -G \\ -H & -A^T \end{bmatrix} \in \mathbf{H}$  with  $\sigma(\mathcal{H}) \cap \text{Im} = \emptyset$ ;  $\epsilon$

Output: the stabilizing solution  $X = X^T \geq 0$  to the CARE (1).

Find an appropriate value  $\hat{\gamma} > 0$ .

Compute  $\widehat{A}_0 \leftarrow I + 2\hat{\gamma}(A_{\hat{\gamma}} + G A_{\hat{\gamma}}^{-T} H)^{-1}$ ,  $\widehat{G}_0 \leftarrow 2\hat{\gamma} A_{\hat{\gamma}}^{-1} G (A_{\hat{\gamma}}^T + H A_{\hat{\gamma}}^{-1} G)^{-1}$ ,

$\widehat{H}_0 \leftarrow 2\hat{\gamma} (A_{\hat{\gamma}}^T + H A_{\hat{\gamma}}^{-1} G)^{-1} H A_{\hat{\gamma}}^{-1}$ ,  $j \leftarrow 0$ ;

Do until convergence:

Compute  $\widehat{A}_{j+1} \leftarrow \widehat{A}_j (I + \widehat{G}_j \widehat{H}_j)^{-1} \widehat{A}_j$ ,  $\widehat{G}_{j+1} \leftarrow \widehat{G}_j + \widehat{A}_j \widehat{G}_j (I + \widehat{H}_j \widehat{G}_j)^{-1} \widehat{A}_j^T$ ,

$\widehat{H}_{j+1} \leftarrow \widehat{H}_j + \widehat{A}_j^T (I + \widehat{H}_j \widehat{G}_j)^{-1} \widehat{H}_j \widehat{A}_j$ ,  $j \leftarrow j + 1$ ;

If  $\|\widehat{H}_j - \widehat{H}_{j-1}\| \leq \epsilon \|\widehat{H}_j\|$ , Stop;

End

Set  $X \leftarrow \widehat{H}_j$ .

## 2.2. Convergence of SDA

Let  $\widehat{\mathcal{N}} = \begin{bmatrix} \widehat{A} & 0 \\ -\widehat{H} & I \end{bmatrix}$ ,  $\widehat{\mathcal{L}} = \begin{bmatrix} I & \widehat{G} \\ 0 & \widehat{A}^T \end{bmatrix}$ , where  $\widehat{G} = \widehat{G}^T$ ,  $\widehat{H} = \widehat{H}^T$ . Suppose  $\widehat{\mathcal{N}} - \lambda \widehat{\mathcal{L}}$  has no eigenvalues on the unit circle and there exist nonsingular  $Q, Z$  such that

$$Q\widehat{\mathcal{N}}Z = \begin{bmatrix} J_s & 0 \\ 0 & I \end{bmatrix}, \quad Q\widehat{\mathcal{L}}Z = \begin{bmatrix} I & 0 \\ 0 & J_s \end{bmatrix}, \quad (13)$$

where the spectrum  $\lambda(J_s) \in O_s \equiv \{\lambda : |\lambda| < 1\}$ . In the following we quote the convergence results for the SDA algorithm from [24,25].

**Theorem 2.1** [24]. Let  $\widehat{\mathcal{N}} = \begin{bmatrix} \widehat{A} & 0 \\ -\widehat{H} & I \end{bmatrix}$  and  $\widehat{\mathcal{L}} = \begin{bmatrix} I & \widehat{G} \\ 0 & \widehat{A}^T \end{bmatrix}$ , where  $\widehat{G} = \widehat{G}^T$ ,  $\widehat{H} = \widehat{H}^T$ . Suppose  $\widehat{\mathcal{N}} - \lambda \widehat{\mathcal{L}}$  has no eigenvalues on the unit circle and there exist nonsingular  $Q, Z$  such that (13) holds. Denote  $Z = \begin{bmatrix} Z_1 & Z_3 \\ Z_2 & Z_4 \end{bmatrix}$ ,  $Z_i \in \mathbb{R}^{n \times n}$  for  $i = 1, 2, 3, 4$ . If  $Z_1$  and  $Z_4$  are invertible, then the sequences  $\{\widehat{A}_j, \widehat{H}_j, \widehat{G}_j\}$  computed by the SDA algorithm satisfy

- (i)  $\|\widehat{A}_j\| = O(\|J_s^{2^j}\|) \rightarrow 0$  as  $j \rightarrow \infty$ ,
- (ii)  $\widehat{H}_j \rightarrow X$ , where  $X$  solves the DARE (12) :

$$X = \widehat{A}^T X (I + \widehat{G}X)^{-1} \widehat{A} + \widehat{H},$$

- (iii)  $\widehat{G}_j \rightarrow Y$ , where  $Y$  solves the dual DARE

$$Y = \widehat{A}Y(I + \widehat{H}Y)^{-1} \widehat{A}^T + \widehat{G}. \quad (14)$$

Moreover, the convergence rate in (i)–(iii) above is  $O(|\lambda_n|^{2^j})$ , where  $|\lambda_1| \leq \dots \leq |\lambda_n| < 1 < |\lambda_n|^{-1} \leq \dots \leq |\lambda_1|^{-1}$  with  $\lambda_i, \lambda_i^{-1}$  being the eigenvalues of  $\widehat{\mathcal{N}} - \lambda \widehat{\mathcal{L}}$  (including 0 and  $\infty$ ).

The following lemma proves that the stabilizability and detectability properties are preserved by the SDA throughout its iterative process. From Lemma A.1 and the SDA algorithm, the matrices  $\widehat{G}_j$  and  $\widehat{H}_j$  are positive semi-definite for each  $j \geq 1$ .

**Lemma 2.2** [24]. The stabilizability of  $(\widehat{A}, \widehat{B})$  implies that  $(\widehat{A}_j, \widehat{B}_j)$  is stabilizable, where  $\widehat{G}_j = \widehat{B}_j \widehat{B}_j^T \geq 0$  is a FRD of  $\widehat{G}_j$  for each  $j \geq 1$ . The detectability of  $(\widehat{A}, \widehat{C})$  implies that  $(\widehat{A}_j, \widehat{C}_j)$  is detectable, where  $\widehat{H}_j = \widehat{C}_j^T \widehat{C}_j \geq 0$  is a FRD of  $\widehat{H}_j$  for each  $j \geq 1$ .

**Theorem 2.3** [24]. Let  $\widehat{\mathcal{N}} = \begin{bmatrix} \widehat{A} & 0 \\ -\widehat{H} & I \end{bmatrix}$  and  $\widehat{\mathcal{L}} = \begin{bmatrix} I & \widehat{G} \\ 0 & \widehat{A}^T \end{bmatrix}$ , where the matrices  $\widehat{G} = \widehat{B} \widehat{B}^T \geq 0$  (FRD) and  $\widehat{H} = \widehat{C}^T \widehat{C} \geq 0$  (FRD). Assume that  $(\widehat{A}, \widehat{B})$  is stabilizable



and  $(\widehat{A}, \widehat{C})$  is detectable. Then the sequences  $\{\widehat{A}_j, \widehat{H}_j, \widehat{G}_j\}$  computed by the SDA satisfy (i), (ii), (iii) as in Theorem 2.1.

**Remarks.** Theorem 2.1 directly proves, under the assumptions that  $\widehat{\mathcal{N}} - \lambda \widehat{\mathcal{L}}$  have no unit modulo eigenvalues and  $Z_1, Z_4$  are invertible, that the sequences  $\{\widehat{A}_j, \widehat{H}_j, \widehat{G}_j\}$  generated by the SDA converge to zero and the unique s.p.s.d. solutions of the DAREs in (12) and (14), respectively. Lemma 2.2 shows the preservation of stabilizability and detectability of the iterates  $(\widehat{A}_j, \widehat{G}_j, \widehat{H}_j)$  generated by the SDA. Furthermore, in Theorem 2.3, we see that the assumptions in Theorem 2.1 are weaker than conditions (S) and (D). This distinction of preserving the symplectic structure in SSF, as well as the difference in operation counts, are responsible for the superior performance of the SDA.

On the other hand, for a given  $\mathcal{H} = \begin{bmatrix} A & -G \\ -H & -A^T \end{bmatrix} \in \mathbf{H}$  with  $\sigma(\mathcal{H}) \cap \text{Im} = \emptyset$ , the matrix sign function of  $\mathcal{H}$  can also be used to develop a structure-preserving method for computing the stabilizing solution of CARE (1). A thorough discussion and the details of practical implementation are given in [21,41]. The main MSGM algorithm is described as follows. Other modified versions can be found in [5,8,9,22,29] and references therein.

**Matrix sign function algorithm:** [7,11–14,20,21,27,33,46]

Input:  $\mathcal{H} = \begin{bmatrix} A & -G \\ -H & -A^T \end{bmatrix} \in \mathbf{H}$  with  $\sigma(\mathcal{H}) \cap \text{Im} = \emptyset$ ;  $\epsilon$ .

Output: the stabilizing solution  $X = X^T \geq 0$  to the CARE (1).

Let  $\mathcal{H}_0 \leftarrow \mathcal{H}$ ,  $j \leftarrow 0$ .

Do until convergence:

    Compute  $\mathcal{H}_{j+1} \leftarrow \frac{1}{2}(\mathcal{H}_j + \mathcal{H}_j^{-1})$ ,  $j \leftarrow j + 1$ ;

    If  $\|\mathcal{H}_j - \mathcal{H}_{j-1}\| \leq \epsilon \|\mathcal{H}_j\|$ , Stop;

End

$\text{sgn}(\mathcal{H}) \leftarrow \mathcal{H}_j$ ;

Solve  $(I - \text{sgn}(\mathcal{H})) \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} = 0$ ;

Compute  $X \leftarrow X_2 X_1^{-1}$ .

#### Remarks

- (i) Notice that  $\begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$  spans the stable invariant subspace of the  $\mathcal{H}$ .
- (ii) Both the SDA and the matrix sign function algorithm require  $\frac{32}{3}n^3$  flops for each iterative step.
- (iii) When working with the Hamiltonian matrix  $\mathcal{H}$ , a more efficient and structure-preserving version of the classical matrix sign function iteration can be derived

by working only with symmetric matrices  $J\mathcal{H}$ . Details may be consulted in [21,41].

### 3. Practical implementation of SDA

#### 3.1. Selection of $\hat{\gamma}$

Here we first derived the forward error bounds of matrices  $\hat{A}_0 \equiv \hat{A}$ ,  $\hat{G}_0 \equiv \hat{G}$  and  $\hat{H}_0 \equiv \hat{H}$  given in (9)–(11), respectively. According to these forward errors, we can design a numerical scheme to determine an appropriate value  $\hat{\gamma} > 0$ . In the following roundoff analysis, we use  $\text{fl}(\cdot)$  to denote computed floating point values. The quantity  $\mathbf{u}$  is the *unit roundoff* (or machine precision), which is typically of order  $10^{-8}$  or  $10^{-16}$  in single and double precision computer arithmetic, respectively. When  $A$  and  $B$  are  $m \times n$  real matrices, the matrix  $B := |A|$  if  $b_{ij} = |a_{ij}|$  for all  $i, j$ , and  $A \leq B$  if  $a_{ij} \leq b_{ij}$  for all  $i, j$ . The 1-,  $\infty$ - and Frobenius matrix norms are denoted by  $\|\cdot\|_1$ ,  $\|\cdot\|_\infty$  and  $\|\cdot\|_F$ , respectively.

We assume that the LU factorizations of  $A_\gamma$  and  $W_\gamma \equiv A_\gamma + GA_\gamma^{-T}H$  are computed by Gaussian elimination with partial pivoting (GEPP). We write these computed LU factors as  $L_A, U_A, L_W$  and  $U_W$ , respectively. Recall that

$$A_\gamma + \Delta A_\gamma = L_A U_A, \quad |\Delta A_\gamma| \leq \gamma_n |L_A| |U_A|, \quad (15)$$

$$W_\gamma + \Delta W_\gamma = L_W U_W, \quad |\Delta W_\gamma| \leq \gamma_n |L_W| |U_W| \quad (16)$$

with  $\gamma_n := n\mathbf{u}/(1 - n\mathbf{u})$  (see, e.g., [32–Theorem 9.3]). Then we have

$$\text{fl}(W_\gamma^{-1}) = W_\gamma^{-1} + E_1, \quad |E_1| \leq c_n \mathbf{u} |W_\gamma^{-1}| |L_W| |U_W| |\text{fl}(W_\gamma^{-1})|, \quad (17)$$

where  $c_n$  is a modest constant. From (17), the forward error bound in evaluating  $\hat{A}$  in (9) is

$$\begin{aligned} \text{fl}(\hat{A}) &= \hat{A} + E_2, \\ |E_2| &\leq 4\gamma c_n \mathbf{u} |W_\gamma^{-1}| |L_W| |U_W| |\text{fl}(W_\gamma^{-1})| + \mathbf{u} |\hat{A}| + O(\mathbf{u}^2). \end{aligned} \quad (18)$$

Furthermore, from (15), we have

$$\begin{aligned} \hat{G}_\gamma &\equiv \text{fl}(2\gamma A_\gamma^{-1}G) = 2\gamma A_\gamma^{-1}G + E_3, \\ |E_3| &\leq 2\gamma c_n \mathbf{u} |A_\gamma^{-1}| |L_A| |U_A| |\hat{G}_\gamma|, \end{aligned} \quad (19)$$

hence the forward error bound in evaluating  $\hat{G}$  in (10) is

$$\begin{aligned} \text{fl}(\hat{G}) &= \hat{G} + E_4, \\ |E_4| &\leq 2\gamma c_n \mathbf{u} |A_\gamma^{-1}| |L_A| |U_A| |\hat{G}_\gamma| |W_\gamma^{-1}|^T \\ &\quad + c_n \mathbf{u} |\text{fl}(\hat{G})| |U_W^T| |L_W^T| |W_\gamma^{-1}|^T. \end{aligned} \quad (20)$$

Finally, from (15), we have

$$\begin{aligned}\widehat{H}_\gamma &\equiv \text{fl}(2\gamma H A_\gamma^{-1}) = 2\gamma H A_\gamma^{-1} + E_5, \\ |E_5| &\leq 2\gamma c_n \mathbf{u} |\widehat{H}_\gamma| |L_A| |U_A| |A_\gamma^{-1}| \end{aligned} \quad (21)$$

and the forward error bound in evaluating  $\widehat{H}$  in (11) is

$$\begin{aligned}\text{fl}(\widehat{H}) &= \widehat{H} + E_6, \\ |E_6| &\leq 2\gamma c_n \mathbf{u} |W_\gamma^{-1}|^T |\widehat{H}_\gamma| |L_A| |U_A| |A_\gamma^{-1}| \\ &\quad + c_n \mathbf{u} |W_\gamma^{-1}|^T |U_W^T| |L_W^T| |\text{fl}(\widehat{H})|. \end{aligned} \quad (22)$$

For GEPP, we have in practice  $\| |L_A| |U_A| \|_\infty \approx \|A_\gamma\|_\infty$  and  $\| |L_W| |U_W| \|_\infty \approx \|W_\gamma\|_\infty$ , and it follows from (18), (20) and (22) that:

$$\|\text{fl}(\widehat{A}) - \widehat{A}\|_\infty \lesssim 4c_n \mathbf{u} \gamma \kappa_\infty(W_\gamma) \|\text{fl}(W_\gamma^{-1})\|_\infty + \mathbf{u} \|\widehat{A}\|_\infty + O(\mathbf{u}^2), \quad (23)$$

$$\begin{aligned}\|\text{fl}(\widehat{G}) - \widehat{G}\|_\infty &\lesssim 2c_n \mathbf{u} \gamma \kappa_\infty(A_\gamma) \|W_\gamma^{-1}\|_1 \|\widehat{G}_\gamma\|_\infty \\ &\quad + c_n \mathbf{u} \kappa_1(W_\gamma) \|\text{fl}(\widehat{G})\|_\infty, \end{aligned} \quad (24)$$

$$\begin{aligned}\|\text{fl}(\widehat{H}) - \widehat{H}\|_\infty &\leq 2c_n \mathbf{u} \gamma \kappa_\infty(A_\gamma) \|W_\gamma^{-1}\|_1 \|\widehat{H}_\gamma\|_\infty \\ &\quad + c_n \mathbf{u} \kappa_1(W_\gamma) \|\text{fl}(\widehat{H})\|_\infty, \end{aligned} \quad (25)$$

where  $\kappa_1(W_\gamma) \equiv \|W_\gamma\|_1 \|W_\gamma^{-1}\|_1$ ,  $\kappa_\infty(W_\gamma) \equiv \|W_\gamma\|_\infty \|W_\gamma^{-1}\|_\infty$  and  $\kappa_\infty(A_\gamma) \equiv \|A_\gamma\|_\infty \|A_\gamma^{-1}\|_\infty$ .

In order to control the forward error bounds of  $\widehat{A}$ ,  $\widehat{G}$  and  $\widehat{H}$ , we consider the following min–max optimization problem, to determine an optimal value  $\hat{\gamma} > 0$ :

$$\min_{\gamma > 0} F(\gamma) \equiv \max_{i=1,2,3} \{f_i(\gamma)\}, \quad (26)$$

where  $f_1(\gamma) := \gamma \kappa_\infty(W_\gamma)$ ,  $f_2(\gamma) := \gamma \kappa_\infty(A_\gamma)$  and  $f_3(\gamma) := \kappa_1(W_\gamma)$ . Since the condition numbers  $\kappa_\infty(W_\gamma)$ ,  $\kappa_\infty(A_\gamma)$  and  $\kappa_1(W_\gamma)$  approach 1 as  $\gamma \rightarrow \infty$ , it follows that  $F(\gamma)$  becomes unbounded as  $\gamma \rightarrow \infty$ . Extensive numerical experiments on randomly generated matrices indicate that  $F(\gamma)$  is a strictly convex function in the neighborhood of the optimal  $\hat{\gamma}$  where the global minimum of  $F(\gamma)$  occurs. For illustration, we report a sample of graphs of  $f_1(\gamma)$ ,  $f_2(\gamma)$ ,  $f_3(\gamma)$  and  $F(\gamma)$  in Figs. 1 and 2. From Theorem 2.1, we know that if  $\gamma$  approaches 0 and  $\infty$ , the symplectic matrix pair  $(\widehat{\mathcal{N}}, \widehat{\mathcal{L}})$  has eigenvalues close to one, leading to very slow convergence of the SDA. This can be avoided through the min-max optimization problem (26).

We can apply the Fibonacci search method to compute an approximate value of  $\hat{\gamma}$ , see, e.g., [10–p. 272]. Our experience indicates that three to five iterations of Fibonacci search are adequate to obtain a suboptimal yet acceptable approximation to  $\hat{\gamma}$ .

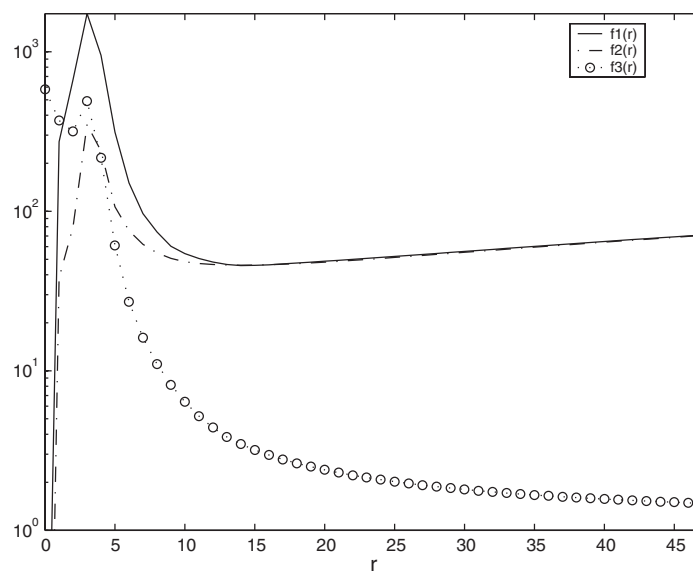


Fig. 1. The graphs of functions  $f_1$ ,  $f_2$  and  $f_3$ .

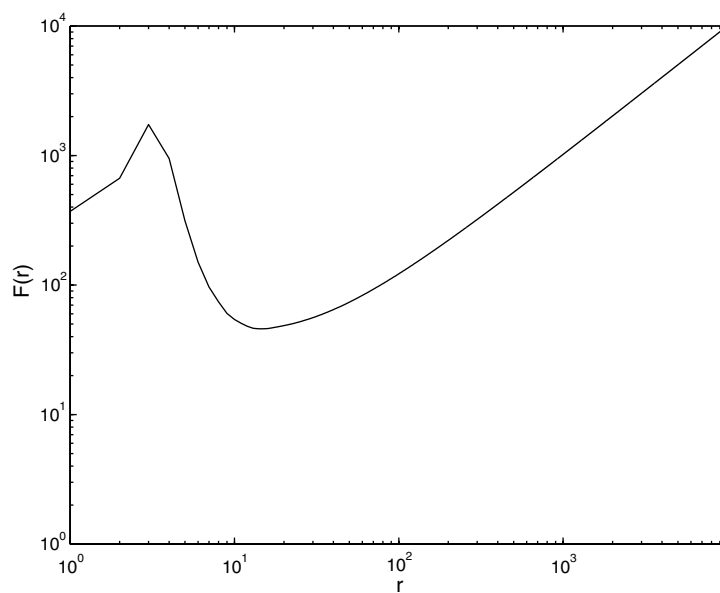


Fig. 2. The graph of  $F(r)$ .

### 3.2. Symmetry of $\widehat{G}_0$ and $\widehat{H}_0$

If the matrices  $G$  and  $H$  are of low-rank, say  $G = gg^T \geq 0$  and  $H = hh^T \geq 0$ , then so are  $\widehat{G}_0$  and  $\widehat{H}_0$ . Indeed, by using the Sherman–Morrison–Woodbury formula (SMWF) twice, it can be seen that

$$\begin{aligned}
 \widehat{G}_0 &= 2\gamma A_\gamma^{-1} G (A_\gamma^T + H A_\gamma^{-1} G)^{-1} \\
 &= 2\gamma A_\gamma^{-1} g g^T (A_\gamma^T + h h^T A_\gamma^{-1} g g^T)^{-1} \\
 &= 2\gamma [A_\gamma^{-1} g g^T (A_\gamma^{-T} - A_\gamma^{-T} h h^T (I + A_\gamma^{-1} g g^T A_\gamma^{-T} h h^T)^{-1} \\
 &\quad \times A_\gamma^{-1} g g^T A_\gamma^{-T})] \\
 &= 2\gamma \{ (A_\gamma^{-1} g) [I - (A_\gamma^{-1} g)^T h h^T (I + (A_\gamma^{-1} g) (A_\gamma^{-1} g)^T h h^T)^{-1} \\
 &\quad \times (A_\gamma^{-1} g)] (A_\gamma^{-1} g)^T \} \\
 &= 2\gamma \{ (A_\gamma^{-1} g) (I + (A_\gamma^{-1} g)^T h h^T (A_\gamma^{-1} g))^{-1} (A_\gamma^{-1} g)^T \} \\
 &= 2\gamma \{ (A_\gamma^{-1} g) (K_g^T K_g)^{-1} (A_\gamma^{-1} g)^T \} \quad (\text{Cholesky decomposition}) \\
 &= 2\gamma (A_\gamma^{-1} g K_g^{-1}) (A_\gamma^{-1} g K_g^{-1})^T.
 \end{aligned}$$

Similarly, by applying the same techniques, we also have

$$\begin{aligned}
 \widehat{H}_0 &= 2\gamma (A_\gamma^T + H A_\gamma^{-1} G)^{-1} H A_\gamma^{-1} \\
 &= 2\gamma (A_\gamma^T + h^T h A_\gamma^{-1} g g^T)^{-1} h^T h A_\gamma^{-1} \\
 &= 2\gamma [A_\gamma^{-T} - A_\gamma^{-T} h^T h A_\gamma^{-1} (I + g g^T A_\gamma^{-T} h^T h A_\gamma^{-1})^{-1} g g^T A_\gamma^{-T}] h^T h A_\gamma^{-1} \\
 &= 2\gamma \{ (h A_\gamma^{-1})^T [I - (h A_\gamma^{-1}) (I + g g^T A_\gamma^{-T} h^T h A_\gamma^{-1})^{-1} \\
 &\quad \times g g^T (h A_\gamma^{-1})^T] (h A_\gamma^{-1}) \} \\
 &= 2\gamma \{ (h A_\gamma^{-1})^T (I + (h A_\gamma^{-1}) g g^T (h A_\gamma^{-1})^T)^{-1} (h A_\gamma^{-1}) \} \\
 &= 2\gamma \{ (h A_\gamma^{-1})^T (K_h K_h^T)^{-1} (h A_\gamma^{-1}) \} \quad (\text{Cholesky decomposition}) \\
 &= 2\gamma (K_h^{-1} h A_\gamma^{-1})^T (K_h^{-1} h A_\gamma^{-1}).
 \end{aligned}$$

### 3.3. Computation of $\widehat{A}_j$ , $\widehat{G}_j$ and $\widehat{H}_j$

We now propose a structured and efficient procedure for the computation of  $\widehat{A}_j$ ,  $\widehat{G}_j$  and  $\widehat{H}_j$  in the SDA algorithm, respectively, where  $\widehat{G}_0 = \widehat{B}_0 \widehat{B}_0^T \geq 0$ ,  $\widehat{H}_0 = \widehat{C}_0^T \widehat{C}_0 \geq 0$  are FRDs. For  $j = 0, 1, 2, \dots$ , we let  $W_j \equiv (I + \widehat{G}_j \widehat{H}_j)^{-1}$ . It is easily seen that  $\widehat{H}_j W_j = W_j^T \widehat{H}_j$  and  $\widehat{G}_j W_j^T = W_j \widehat{G}_j$  are s.p.s.d. for each  $j \geq 1$ . By the SMWF we can derive the formulae

$$W_j = (I + \widehat{G}_j \widehat{H}_j)^{-1} = I - \widehat{B}_j (I + \widehat{B}_j^T \widehat{H}_j \widehat{B}_j)^{-1} \widehat{B}_j^T \widehat{H}_j, \quad (27)$$

$$\widehat{G}_j W_j^T = \widehat{G}_j - \widehat{G}_j \widehat{C}_j^T (I + \widehat{C}_j \widehat{G}_j \widehat{C}_j^T)^{-1} \widehat{C}_j \widehat{G}_j = \widehat{B}_j (I + \widehat{B}_j^T \widehat{H}_j \widehat{B}_j)^{-1} \widehat{B}_j^T, \quad (28)$$

$$W_j^T \widehat{H}_j = \widehat{H}_j - \widehat{H}_j \widehat{B}_j (I + \widehat{B}_j^T \widehat{H}_j \widehat{B}_j)^{-1} \widehat{B}_j^T \widehat{H}_j = \widehat{C}_j^T (I + \widehat{C}_j \widehat{G}_j \widehat{C}_j^T)^{-1} \widehat{C}_j. \quad (29)$$

When the matrices  $B$  and  $C$  start with low ranks in (1), we can improve the efficiency of our computation further by the following compression process. Compute the Cholesky decomposition of the s.p.d. matrices  $W_{G,j} \equiv (I + \widehat{B}_j^T \widehat{H}_j \widehat{B}_j) = K_{B,j}^T K_{B,j}$  and  $W_{H,j} \equiv (I + \widehat{C}_j \widehat{G}_j \widehat{C}_j^T) = K_{C,j} K_{C,j}^T$ , respectively. For  $j = 0, 1, 2, \dots$ , application of (27)–(29) leads to

$$\widehat{A}_{j+1} = \widehat{A}_j^2 - \widehat{A}_j \widehat{B}_j (I + \widehat{B}_j^T \widehat{H}_j \widehat{B}_j)^{-1} \widehat{B}_j^T \widehat{H}_j \widehat{A}_j, \quad (30)$$

$$\begin{aligned} \widehat{G}_{j+1} &= \widehat{G}_j + \widehat{A}_j \widehat{B}_j (I + \widehat{B}_j^T \widehat{H}_j \widehat{B}_j)^{-1} \widehat{B}_j^T \widehat{A}_j^T \\ &= [\widehat{B}_j, \quad \widehat{A}_j \widehat{B}_j K_{B,j}^{-1}] \begin{bmatrix} \widehat{B}_j^T \\ K_{B,j}^{-T} \widehat{B}_j^T \widehat{A}_j^T \end{bmatrix} \equiv \widehat{B}_{j+1} \widehat{B}_{j+1}^T \geq 0 \quad (\text{FRD}) \end{aligned} \quad (31)$$

and

$$\begin{aligned} \widehat{H}_{j+1} &= \widehat{H}_j + \widehat{A}_j^T \widehat{C}_j^T (I + \widehat{C}_j \widehat{G}_j \widehat{C}_j^T)^{-1} \widehat{C}_j \widehat{A}_j \\ &= [\widehat{C}_j^T, \quad \widehat{A}_j^T \widehat{C}_j^T K_{C,j}^{-T}] \begin{bmatrix} \widehat{C}_j \\ K_{C,j}^{-1} \widehat{C}_j \widehat{A}_j \end{bmatrix} \equiv \widehat{C}_{j+1}^T \widehat{C}_{j+1} \geq 0 \quad (\text{FRD}), \end{aligned} \quad (32)$$

where  $\widehat{B}_{j+1}$  and  $\widehat{C}_{j+1}^T$  are the full column rank compressions of  $[\widehat{B}_j, \widehat{A}_j \widehat{B}_j K_{B,j}^{-1}]$  and  $[\widehat{C}_j^T, \widehat{A}_j^T \widehat{C}_j^T K_{C,j}^{-T}]$ , respectively. In general,  $\text{rank}(\widehat{B}_{j+1}) > \text{rank}(\widehat{B}_j)$  and  $\text{rank}(\widehat{C}_{j+1}) > \text{rank}(\widehat{C}_j)$ , and the compression process becomes unprofitable when the ranks of  $\widehat{B}_{j+1}$  and  $\widehat{C}_{j+1}$  approach  $n$ .

### 3.4. Error analysis of SDA

We consider the forward error bounds of the computed matrices  $\widehat{A}_{j+1}$ ,  $\widehat{G}_{j+1}$  and  $\widehat{H}_{j+1}$  in the SDA algorithm for one iterative step  $j$ . Since  $K_{B,j}$  and  $K_{C,j}$  are the computed Cholesky factors of matrices  $W_{G,j}$  and  $W_{H,j}$ , respectively, it follows that:

$$\begin{aligned} \widehat{K}_B &\equiv \text{fl}(K_{B,j}^{-T} \widehat{B}_j^T) = K_{B,j}^{-T} \widehat{B}_j^T + \Delta E_1, \\ |\Delta E_1| &\leq c_1 \mathbf{u} |K_{B,j}^{-T}| |K_{B,j}^T| |\widehat{K}_B| \end{aligned} \quad (33)$$

and

$$\begin{aligned} \widehat{K}_C &\equiv \text{fl}(K_{C,j}^{-1} \widehat{C}_j) = K_{C,j}^{-1} \widehat{C}_j + \Delta \widetilde{E}_1, \\ |\Delta \widetilde{E}_1| &\leq \widetilde{c}_1 \mathbf{u} |K_{C,j}^{-1}| |K_{C,j}| |\widehat{K}_C|, \end{aligned} \quad (34)$$

where  $c_1$  and  $\tilde{c}_1$  are modest constants. Therefore, we have

$$\begin{aligned} \mathfrak{fl}(K_{B,j}^{-T} \widehat{B}_j^T \widehat{A}_j^T) &= \mathfrak{fl}(\widehat{K}_B \widehat{A}_j^T) = K_{B,j}^{-T} \widehat{B}_j^T \widehat{A}_j^T + \Delta E_2, \\ |\Delta E_2| &\leq c_2 \mathbf{u} |K_{B,j}^{-T}| |K_{B,j}^T| |\widehat{K}_B| |\widehat{A}_j^T| \end{aligned} \quad (35)$$

and

$$\begin{aligned} \mathfrak{fl}(K_{C,j}^{-1} \widehat{C}_j \widehat{A}_j) &= \mathfrak{fl}(\widehat{K}_C \widehat{A}_j) = K_{C,j}^{-1} \widehat{C}_j \widehat{A}_j + \Delta \tilde{E}_2, \\ |\Delta \tilde{E}_2| &\leq \tilde{c}_2 \mathbf{u} |K_{C,j}^{-1}| |K_{C,j}| |\widehat{K}_C| |\widehat{A}_j|, \end{aligned} \quad (36)$$

where  $c_2$  and  $\tilde{c}_2$  are modest constants.

If  $\text{rank}(\widehat{B}_j^T) = \ell$ , then from Theorem 18.4 of [32] and (31), there exist an orthogonal matrix  $Q_B \in \mathbb{R}^{2\ell \times 2\ell}$  and a computed upper triangular matrix  $\mathfrak{fl}(\widehat{B}_{j+1}^T)$  with full row rank, such that

$$\begin{bmatrix} \widehat{B}_j^T \\ \mathfrak{fl}(\widehat{K}_B \widehat{A}_j^T) \end{bmatrix} + \begin{bmatrix} \Delta B_1 \\ \Delta B_2 \end{bmatrix} = Q_B \begin{bmatrix} \mathfrak{fl}(\widehat{B}_{j+1}^T) \\ 0 \end{bmatrix}, \quad (37)$$

where  $|\Delta B_j| \leq c_3 \mathbf{u} G_\ell(|\widehat{B}_j^T| + |K_{B,j}^{-T}| |\widehat{B}_j^T| |\widehat{A}_j^T|)$  for  $j = 1, 2$ , with  $c_3$  being a modest constant and  $\|G_\ell\|_F = \frac{1}{2}$ .

From (35) and (37), we have

$$\begin{bmatrix} \widehat{B}_j^T \\ K_{B,j}^{-T} \widehat{B}_j^T \widehat{A}_j^T \end{bmatrix} + \begin{bmatrix} \Delta B_1 \\ \Delta \tilde{B}_2 \end{bmatrix} = Q_B \begin{bmatrix} \mathfrak{fl}(\widehat{B}_{j+1}^T) \\ 0 \end{bmatrix}, \quad (38)$$

where  $|\Delta \tilde{B}_2| \leq c_2 \mathbf{u} |K_{B,j}^{-T}| |K_{B,j}^T| |\widehat{K}_B| |\widehat{A}_j^T| + c_3 \mathbf{u} G_\ell(|\widehat{B}_j^T| + |K_{B,j}^{-T}| |\widehat{B}_j^T| |\widehat{A}_j^T|)$ . Pre-multiplying both sides of (38) by  $Q_B^T$ , it follows that:

$$\begin{bmatrix} \widehat{B}_{j+1}^T \\ 0 \end{bmatrix} + Q_B^T \begin{bmatrix} \Delta B_1 \\ \Delta \tilde{B}_2 \end{bmatrix} = \begin{bmatrix} \mathfrak{fl}(\widehat{B}_{j+1}^T) \\ 0 \end{bmatrix} \quad (39)$$

and we deduce that

$$\|\mathfrak{fl}(\widehat{B}_{j+1}^T) - \widehat{B}_{j+1}^T\|_F \leq c_4 \mathbf{u} \|\widehat{B}_j\|_F + c_5 \mathbf{u} \kappa_s(K_{B,j}^T) \|\widehat{K}_B\|_F \|\widehat{A}_j\|_F, \quad (40)$$

where  $c_4$  and  $c_5$  are modest constants, and  $\kappa_s(K_{B,j}^T) \equiv \| |K_{B,j}^{-T}| |K_{B,j}^T| \|_\infty$  is the Skeel condition number of  $K_{B,j}^T$ . Furthermore, applying a similar argument with the help of (36), we can derive that

$$\|\mathfrak{fl}(\widehat{C}_{j+1}) - \widehat{C}_{j+1}\|_F \leq c_6 \mathbf{u} \|\widehat{C}_j\|_F + c_7 \mathbf{u} \kappa_s(K_{C,j}) \|\widehat{K}_C\|_F \|\widehat{A}_j\|_F, \quad (41)$$

where  $c_6$  and  $c_7$  are modest constants.

On the other hand, it follows from (33) that:

$$\begin{aligned} \mathfrak{fl}(K_{B,j}^{-T} \widehat{B}_j^T \widehat{H}_j \widehat{A}_j) &= K_{B,j}^{-T} \widehat{B}_j^T \widehat{H}_j \widehat{A}_j + \Delta E_3, \\ |\Delta E_3| &\leq c_8 \mathbf{u} |K_{B,j}^{-T}| |K_{B,j}^T| |\widehat{K}_B| |\widehat{H}_j| |\widehat{A}_j|, \end{aligned} \quad (42)$$

where  $c_8$  is a modest constant. From (35) and (42), the forward error bound of computing  $\widehat{A}_{j+1}$  is

$$\begin{aligned} & \|\text{fl}(\widehat{A}_{j+1}) - \widehat{A}_{j+1}\|_F \\ & \leq c_9 \mathbf{u} \|\widehat{A}_j\|_F^2 + c_{10} \mathbf{u} \kappa_s(K_{B,j}^T) \|\widehat{B}_j\|_F^2 \|\widehat{K}_{B,j}^{-1}\|_F^2 \|\widehat{H}_j\|_F \|\widehat{A}_j\|_F^2, \end{aligned} \quad (43)$$

where  $c_9$  and  $c_{10}$  are modest constants.

When the Skeel condition numbers  $\kappa_s(K_{B,j}^T)$  and  $\kappa_s(K_{C,j})$  in (40) and (41) are bounded from above by acceptable numbers, the accumulation of error will be dampened by the fast rate of convergence at the final stage of the iterative process. Danger, if any, lies in the early stage of the process before the  $\lambda_n^{2^j}$  convergence factor dominates. It is unlikely to have any ill-effect, as the accumulated error in the matrix additions and multiplications should be of magnitude around a small multiple of the machine accuracy.

As the SSF properties are preserved in the SDA, any error will be a structured one, only pushing the iteration towards a solution of a neighboring SSF system. Thus the algorithm is stable in this sense, when the errors are not too large and when stabilizability and detectability are maintained. For large  $j$ s, as  $\widehat{A}_j \rightarrow 0$ ,  $\widehat{G}_j$  and  $\widehat{H}_j$  converge to the unique s.p.s.d. solutions of (14) and (12), respectively. Danger again will only come at the initial stage of the iteration. Corresponding checks may be prudent in the algorithm.

#### 4. SDA\_m

A matrix  $A$  is persymmetric when  $A$  is symmetric with respect to the main anti-diagonal [30–p. 193]. When the DARE transformed from the CARE (1) has the additional property that the initial data  $\widehat{A}_0, \widehat{G}_0 = \widehat{H}_0 \in \mathbb{R}^{2\ell \times 2\ell}$  are symmetric and persymmetric, the additional structure can be preserved in a modified version of the SDA (SDA\_m). For simplicity, we consider only when  $\gamma = 1$ . This doubly symmetric type of DAREs appear in the Examples 10 and 17 of Section 5 (originally from [15]).

For convenience, in the SDA, we denote for  $j = 1, 2, \dots$

$$\begin{aligned} A &\equiv \widehat{A}_j, & G &\equiv \widehat{G}_j = \widehat{H}_j, \\ A_+ &\equiv \widehat{A}_{j+1}, & G_+ &\equiv \widehat{G}_{j+1} = \widehat{H}_{j+1}. \end{aligned} \quad (44)$$

Since  $A, G = H$  are symmetric and persymmetric of even order, we write

$$A = \begin{bmatrix} a_1 & a_2 \zeta \\ \zeta a_2 & \zeta a_1 \zeta \end{bmatrix}, \quad G = \begin{bmatrix} g_1 & g_2 \zeta \\ \zeta g_2 & \zeta g_1 \zeta \end{bmatrix}, \quad (45)$$

where  $a_1, a_2, g_1$  and  $g_2 \in \mathbb{R}^{\ell \times \ell}$  are symmetric and  $\zeta = [e_\ell, \dots, e_1]$  with  $e_j$  being the  $j$ th column of  $I_\ell$ . In the SDA, we shall show that  $\widehat{A}, \widehat{G}$  and  $\widehat{H}$  are also symmetric and persymmetric with  $\widehat{G} = \widehat{H}$ , with



$$A_+ = A(I + G^2)^{-1}A = \begin{bmatrix} \hat{a}_1 & \hat{a}_2\zeta \\ \zeta\hat{a}_2 & \zeta\hat{a}_1\zeta \end{bmatrix}, \quad (46)$$

$$G_+ = G + AG(I + G^2)^{-1}A^T = \begin{bmatrix} \hat{g}_1 & \hat{g}_2\zeta \\ \zeta\hat{g}_2 & \zeta\hat{g}_1\zeta \end{bmatrix}. \quad (47)$$

Let

$$q_1 \equiv g_1 + g_2, \quad q_2 \equiv g_1 - g_2, \quad \alpha_1 \equiv a_1 + a_2, \quad \alpha_2 \equiv a_1 - a_2. \quad (48)$$

Simple manipulation leads to

$$\hat{a}_1 = \frac{1}{2}[\alpha_1(I + q_1^2)^{-1}\alpha_1 + \alpha_2(I + q_2^2)^{-1}\alpha_2], \quad (49)$$

$$\hat{a}_2 = \frac{1}{2}[\alpha_1(I + q_1^2)^{-1}\alpha_1 - \alpha_2(I + q_2^2)^{-1}\alpha_2], \quad (50)$$

$$\hat{g}_1 = g_1 + \frac{1}{2}[\alpha_1(I + q_1^2)^{-1}q_1\alpha_1 + \alpha_2(I + q_2^2)^{-1}q_2\alpha_2], \quad (51)$$

$$\hat{g}_2 = g_2 + \frac{1}{2}[\alpha_1(I + q_1^2)^{-1}q_1\alpha_1 - \alpha_2(I + q_2^2)^{-1}q_2\alpha_2]. \quad (52)$$

Furthermore let

$$q_1 = [U_1, V_1] \begin{bmatrix} \Sigma_1 & 0 \\ 0 & -\Gamma_1 \end{bmatrix} \begin{bmatrix} U_1^T \\ V_1^T \end{bmatrix}, \quad q_2 = [U_2, V_2] \begin{bmatrix} \Sigma_2 & 0 \\ 0 & -\Gamma_2 \end{bmatrix} \begin{bmatrix} U_2^T \\ V_2^T \end{bmatrix} \quad (53)$$

be the spectral decompositions of  $q_1$  and  $q_2$ , respectively, with  $\Sigma_1, \Gamma_1, \Sigma_2$  and  $\Gamma_2$  being nonnegative diagonal matrices. Then  $\hat{a}_1, \hat{a}_2, \hat{g}_1$  and  $\hat{g}_2$  in (49)–(52) can be computed by the following symmetric forms:

$$\begin{aligned} \xi_1 &\equiv \alpha_1 U_1 (I + \Sigma_1^2)^{-1} U_1^T \alpha_1 - \alpha_1 V_1 (I + \Gamma_1^2)^{-1} V_1^T \alpha_1, \\ \xi_2 &\equiv \alpha_2 U_2 (I + \Sigma_2^2)^{-1} U_2^T \alpha_2 - \alpha_2 V_2 (I + \Gamma_2^2)^{-1} V_2^T \alpha_2, \\ \hat{a}_1 &= \frac{1}{2}\{\xi_1 + \xi_2\}, \quad \hat{a}_2 = \frac{1}{2}\{\xi_1 - \xi_2\}; \end{aligned} \quad (54)$$

$$\begin{aligned} \eta_1 &\equiv \alpha_1 U_1 (I + \Sigma_1^2)^{-1} \Sigma_1 U_1^T \alpha_1 - \alpha_1 V_1 (I + \Gamma_1^2)^{-1} \Gamma_1 V_1^T \alpha_1, \\ \eta_2 &\equiv \alpha_2 U_2 (I + \Sigma_2^2)^{-1} \Sigma_2 U_2^T \alpha_2 - \alpha_2 V_2 (I + \Gamma_2^2)^{-1} \Gamma_2 V_2^T \alpha_2, \\ \hat{g}_1 &= g_1 + \frac{1}{2}\{\eta_1 + \eta_2\}, \quad \hat{g}_2 = g_2 + \frac{1}{2}\{\eta_1 - \eta_2\}. \end{aligned} \quad (55)$$

The SDA\_m computes  $\hat{A}, \hat{G}$  in (46) and (47) using the symmetric forms (54) and (55) and considerably improves the accuracy of Examples 10 and 17 in the next section.

## 5. Numerical examples

For the tables in the following examples, data for various methods are lists in columns with obvious headings. The heading “care” is for the care command in MATLAB [40], “MSGM” is for the matrix sign function method [21], and “SDA” (or “SDA\_m”) stands for our SDA (or SDA\_m) algorithm. There is no iteration numbers to report for care and an ‘\*’ in the tables indicates a failure of convergence in obtaining a solution. In the graphs, “ratio\_care” and “ratio\_MSGM” are the ratio of the CPU-times for care and MSGM to that of the SDA, respectively. For the comparison of residuals, the “normalized” residual (NRes) formula is applied in the numerical examples, i.e.,

$$\text{NRes} \equiv \frac{\|A^T \tilde{X} + \tilde{X} A^T - \tilde{X} G \tilde{X} + H\|}{\|A^T \tilde{X}\| + \|\tilde{X} A^T\| + \|\tilde{X} G \tilde{X}\| + \|H\|},$$

where  $\tilde{X}$  is an approximate solution and  $\|\cdot\|$  denotes the 2-norm for matrices.

Some numerical examples from [15] involved very large data sets, which have not been repeated here. Twenty examples were presented in [23]. We retain the numbering of examples in [23], comment upon all of them but present only five representative ones in this paper.

In the MSGM, the scaling strategy suggested in [21] was implemented. For a fairer comparison, similar convergence criteria were used in all the methods and the solutions were not refined.

All computations were performed using MATLAB/Version 6.0 on a Compaq/DS20 workstation. The machine precision is  $2.22 \times 10^{-16}$ .

**Example 5.** The example is identical to Example 5 of [15], which has been presented originally in [45]. This is a 9th-order continuous state space model of a tubular ammonia reactor. The actual system matrices are

$$A = \begin{bmatrix} -4.019 & 5.12 & 0 & 0 & -2.082 & 0 & 0 & 0 & 0.87 \\ -0.346 & 0.986 & 0 & 0 & -2.34 & 0 & 0 & 0 & 0.97 \\ -7.909 & 15.407 & -4.096 & 0 & -6.45 & 0 & 0 & 0 & 2.68 \\ -21.816 & 35.606 & -0.339 & -3.87 & -17.8 & 0 & 0 & 0 & 7.39 \\ -60.196 & 98.188 & -7.907 & 0.34 & -53.008 & 0 & 0 & 0 & 20.4 \\ 0 & 0 & 0 & 0 & 94.0 & -147.2 & 0 & 53.2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 94.0 & -147.2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 12.8 & 0 & -31.6 & 0 \\ 0 & 0 & 0 & 0 & 12.8 & 0 & 0 & 18.8 & -31.6 \end{bmatrix},$$

$$B^T = \begin{bmatrix} 0.010 & 0.003 & 0.009 & 0.024 & 0.068 & 0 & 0 & 0 & 0 \\ -0.011 & -0.021 & -0.059 & -0.162 & -0.445 & 0 & 0 & 0 & 0 \\ -0.151 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix},$$

$$H = I_9, \quad R = I_3.$$

The numerical results are given in Table 1.

Table 1  
Results for Example 5

	SDA	MSGM	care
NRes	$1.68 \times 10^{-15}$	$1.73 \times 10^{-13}$	$4.64 \times 10^{-14}$
Iter. no.	9	8	–

Table 2  
Results for Example 6

	SDA	MSGM	care
NRes	$5.78 \times 10^{-13}$	$3.11 \times 10^{-8}$	$1.91 \times 10^{-12}$
Iter. no.	10	9	–

**Example 6.** The example is identical to Example 6 of [15], which has been presented originally in [26]. This control problem for a J-100 jet engine is a special case of a multivariable servomechanism problem. To save space, we shall not list the system matrices here. We report the numerical results in Table 2.

**Example 10.** The example is identical to Example 10 of [15], which has been presented originally in [6]. Here, the system matrices are

$$A = \begin{bmatrix} \varepsilon + 1 & 1 \\ 1 & \varepsilon + 1 \end{bmatrix}, \quad G = I_2, \quad H = \begin{bmatrix} \varepsilon^2 & 0 \\ 0 & \varepsilon^2 \end{bmatrix}.$$

The exact stabilizing solution  $X$  is given by

$$\begin{aligned} x_{11} = x_{22} &= \frac{1}{2} \left[ 2(\varepsilon + 1) + \sqrt{2(\varepsilon + 1)^2 + 2} + \sqrt{2\varepsilon} \right], \\ x_{12} = x_{21} &= x_{11} / [x_{11} - (\varepsilon + 1)]. \end{aligned} \quad (56)$$

The corresponding DARE is doubly symmetric and the SDA\_m was applied (see details in Section 4). The numerical results with  $\varepsilon = 1, 10^{-3}, 10^{-5}$  and  $10^{-7}$  are given in Table 3.

**Example 11.** The example is identical to Example 11 of [15], which has been presented originally in [35]. This example represents an algebraic Riccati equation arising from a  $H_\infty$ -control problem [49]. Let

$$A = \begin{bmatrix} 3 - \varepsilon & 1 \\ 4 & 2 - \varepsilon \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad R = 1, \quad H = \begin{bmatrix} 4\varepsilon - 11 & 2\varepsilon - 5 \\ 2\varepsilon - 5 & 2\varepsilon - 2 \end{bmatrix}.$$

The matrix

$$X = \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix}$$

Table 3  
Results for Example 10

		SDA	SDA_m	MSGM	care
$\varepsilon = 1$	NRes	$0.00 \times 10^0$	$0.00 \times 10^0$	$4.69 \times 10^{-16}$	$9.36 \times 10^{-17}$
	Rel. err.	$1.96 \times 10^{-16}$	$1.96 \times 10^{-16}$	$8.80 \times 10^{-16}$	$3.83 \times 10^{-16}$
	Iter. no.	4	4	2	–
$\varepsilon = 10^{-3}$	NRes	$1.58 \times 10^{-14}$	$1.43 \times 10^{-16}$	$1.11 \times 10^{-13}$	$9.20 \times 10^{-17}$
	Rel. err.	$1.82 \times 10^{-11}$	$2.22 \times 10^{-16}$	$2.22 \times 10^{-13}$	$4.08 \times 10^{-16}$
	Iter. no.	16	13	12	–
$\varepsilon = 10^{-5}$	NRes	$2.28 \times 10^{-12}$	$1.11 \times 10^{-16}$	$1.07 \times 10^{-11}$	$5.53 \times 10^{-17}$
	Rel. err.	$7.16 \times 10^{-7}$	$1.76 \times 10^{-16}$	$2.14 \times 10^{-11}$	$2.60 \times 10^{-16}$
	Iter. no.	22	19	18	–
$\varepsilon = 10^{-7}$	NRes	$1.49 \times 10^{-10}$	$1.32 \times 10^{-16}$	$3.31 \times 10^{-9}$	$2.06 \times 10^{-17}$
	Rel. err.	$6.04 \times 10^{-8}$	$4.44 \times 10^{-16}$	$6.63 \times 10^{-9}$	$1.36 \times 10^{-16}$
	Iter. no.	12	26	20	–

Table 4  
Results for Example 11

		SDA	MSGM	care
$\varepsilon = 1$	NRes	$0.00 \times 10^0$	$1.69 \times 10^{-16}$	$1.97 \times 10^{-16}$
	Rel. err.	$1.26 \times 10^{-16}$	$1.25 \times 10^{-15}$	$9.68 \times 10^{-16}$
	Iter. no.	5	2	–
$\varepsilon = 0$	NRes	$3.06 \times 10^{-16}$	*	$5.06 \times 10^{-17}$
	Rel. err.	$2.66 \times 10^{-9}$	*	$7.68 \times 10^{-9}$
	Iter. no.	28	*	–

is the stabilizing solution for  $\varepsilon > 0$ . For  $\varepsilon = 0$ , the solution  $X$  is obtained by an  $H$ -invariant *Lagrangian* subspace, i.e., a solution in the sense of  $H_\infty$ -control. The numerical results with  $\varepsilon = 1, 0$  are given in Table 4.

**Example 12.** The example is identical to Example 12 of [15], which has been presented originally in [34]. Let

$$V = I - \frac{2}{3}vv^T, \quad v^T = [1 \quad 1 \quad 1]; \quad A_0 = \varepsilon \operatorname{diag}(1, 2, 3),$$

$$H_0 = \operatorname{diag}(\varepsilon^{-1}, 1, \varepsilon);$$

we have

$$A = VA_0V, \quad G = \varepsilon^{-1}I_3, \quad H = VH_0V.$$

The solution is

$$X = V \operatorname{diag}(x_1, x_2, x_3)V,$$

Table 5  
Results for Example 12

		SDA	MSGM	care
$\varepsilon = 1$	NRes	$2.01 \times 10^{-16}$	$1.78 \times 10^{-15}$	$3.00 \times 10^{-16}$
	Rel. err.	$4.33 \times 10^{-16}$	$2.78 \times 10^{-15}$	$5.03 \times 10^{-16}$
	Iter. no.	6	4	–
$\varepsilon = 10^6$	NRes	$1.62 \times 10^{-15}$	$2.22 \times 10^{-4}$	$2.19 \times 10^{-15}$
	Rel. err.	$2.58 \times 10^{-15}$	$6.33 \times 10^{-4}$	$4.92 \times 10^{-15}$
	Iter. no.	11	10	–

where

$$x_1 = \varepsilon^2 + \sqrt{\varepsilon^4 + 1}, \quad x_2 = 2\varepsilon^2 + \sqrt{4\varepsilon^4 + \varepsilon}, \quad x_3 = 3\varepsilon^2 + \sqrt{9\varepsilon^4 + \varepsilon^2}.$$

The numerical results with  $\varepsilon = 1, 10^6$  are given in Table 5.

**Example 15.** The example is identical to Example 15 of [15], which has been presented originally in [39–Example 4] and [3]. This example arises from a mathematical model of position and velocity control for a string of high-speed vehicles. If  $N$  vehicles are to be controlled, the size of the system matrices will be  $n = 2N - 1$ , the number of control inputs will be  $m = N$ , and the number of outputs will be  $p = N - 1$ , respectively. The comparison of normalized residuals are reported in Table 6 for  $N = 5, 20, 60, 100, 140$  and 180. Fig. 3 reports the comparison of CPU times for care, MSGN and the SDA.

**Example 17.** The example is identical to Example 17 of [15], which has been presented originally in [39–Example 6]. The system matrices are

$$A = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & 0 \\ 0 & \cdots & \cdots & 0 & 1 \\ 0 & \cdots & \cdots & 0 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 0 \\ \vdots \\ \vdots \\ 0 \\ 1 \end{bmatrix}, \quad R = r, \quad C^T = \sqrt{q} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ \vdots \\ 0 \end{bmatrix}.$$

It is known from [39] that  $x_{1n} = \sqrt{qr}$ . Therefore, we may use the relative error in  $x_{1n}$ , i.e.,  $RE \equiv (|x_{1n} - \sqrt{qr}|)/\sqrt{qr}$ , as an indicator of the accuracy of the results. The corresponding DARE is doubly symmetric and the SDA\_m was applied (see details in Section 4). Table 7 reports the comparison of normalized residuals computed by SDA, SDA\_m and care for  $n = 6, 12, 18, 24, 30$ . We also report the comparison of relative errors in  $x_{1n}$  computed by above three algorithms in Table 8.

Table 6  
Comparison of normalized residuals for Example 15

		SDA	MSGM	care
$N = 5$	NRes	$1.61 \times 10^{-16}$	$8.75 \times 10^{-15}$	$2.53 \times 10^{-15}$
	Iter. no.	5	6	–
$N = 20$	NRes	$3.85 \times 10^{-16}$	$3.55 \times 10^{-14}$	$6.15 \times 10^{-15}$
	Iter. no.	5	6	–
$N = 60$	NRes	$1.53 \times 10^{-15}$	$2.32 \times 10^{-13}$	$8.14 \times 10^{-15}$
	Iter. no.	7	8	–
$N = 100$	NRes	$2.15 \times 10^{-15}$	$6.62 \times 10^{-13}$	$2.55 \times 10^{-14}$
	Iter. no.	8	9	–
$N = 140$	NRes	$3.05 \times 10^{-15}$	$6.50 \times 10^{-12}$	$3.60 \times 10^{-14}$
	Iter. no.	8	9	–
$N = 180$	NRes	$1.25 \times 10^{-14}$	$4.64 \times 10^{-12}$	$2.01 \times 10^{-13}$
	Iter. no.	9	9	–

### 5.1. Comments on numerical results

We have tested 20 examples in [23] to illustrate the accuracy and efficiency of the SDA applied to CAREs, in comparison to the MSGM [21] and care in MATLAB [40]. Some of these examples have parameters to vary their sizes or conditioning. In what follows, we shall comment upon all the examples in [23], thus retaining the old labelling of the examples:

- (1) Comparing with care for all the examples, solutions with better or comparable accuracy were obtained using the SDA in far less time. This comparison has been difficult as care yields no iteration numbers and the CPU time information from MATLAB is not always accurate.
- (2) The best indication of the efficiency of the SDA over care comes from Example 15 (with varying dimension  $n$ ), where care required two to eight times more CPU times than the SDA. This is consistent with the findings in [24] for DAREs. Keep in mind that the SDA requires far less number of flops than care in each iteration, as the operations in the SDA are performed in  $\mathbb{R}^{n \times n}$  whereas those for care are carried out in  $\mathbb{R}^{2n \times 2n}$ .
- (3) For examples with varying conditioning, such as Examples 9–14, 17 and 18, the SDA out-performed care and converges to more accurate solutions in less time. For the ill-conditioned Example 20, care failed while the SDA succeeded without difficulty.
- (4) In Example 11 (in  $H_\infty$  control), some eigenvalues were numerically on the imaginary axis and assumptions in the theory were practically violated. The stronger structure-preserving property of the SDA enabled it to produce an

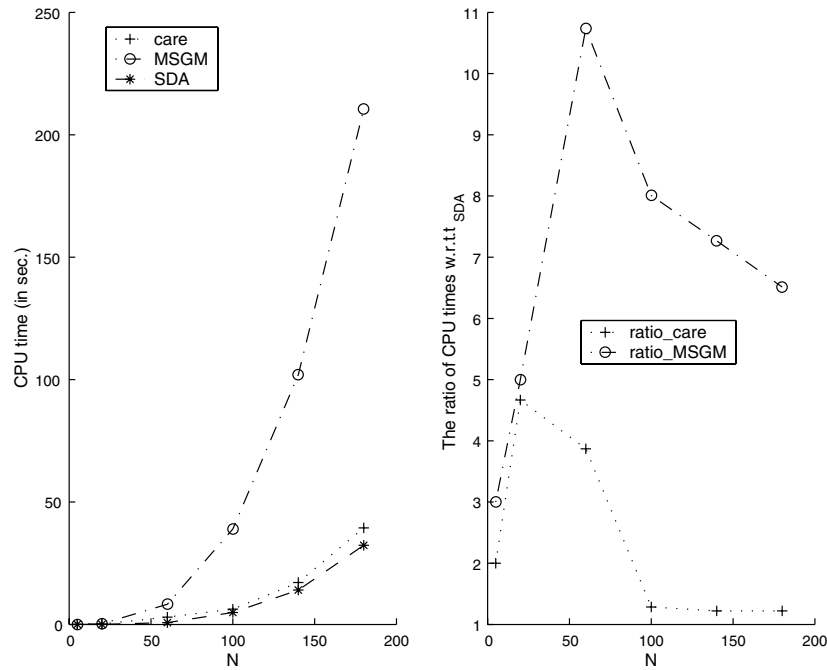


Fig. 3. Comparison of CPU times for Example 15.

accurate solution when the MSGM failed. Somehow, care produced a slightly less accurate solution using much more CPU time.

- (5) In Examples 10 and 17, the CAREs gave rise to DAREs which were “doubly symmetric” (see Section 4 for details). The SDA\_m improved the efficiency of the SDA for these examples, obtaining comparable accuracy for Example 10 while out-performing care for Example 17.
- (6) Comparing to the MSGM for ill-conditioned problems, the SDA performed better in terms of accuracies or number of iterations. This is consistent with the fact that while both the SDA and MSGM are structure-preserving, the former preserves more structure than the latter. For some well-conditioned problems, the efficiency and accuracy of the SDA and MSGM are comparable. For a few simple small examples, the MSGM converged quickly and was superior to the SDA. Note that the work involved in an iteration for either method is similar.
- (7) The MSGM, with similar operations count to SDA, was generally more efficient than care, especially for well-conditioned problems. For ill-conditioned problems (such as Example 10), the MSGM was sometimes less accurate than care.

Table 7  
Comparison of normalized residuals for Example 17

	$n$	NRes_SDA	NRes_SDA_m	NRes_MSGM	NRes_care
$q, r = 1$	6	$4.50 \times 10^{-15}$	$3.56 \times 10^{-16}$	$8.87 \times 10^{-15}$	$1.80 \times 10^{-14}$
	12	$3.63 \times 10^{-10}$	$3.22 \times 10^{-14}$	$9.68 \times 10^{-12}$	$1.23 \times 10^{-11}$
	18	$9.47 \times 10^{-5}$	$1.83 \times 10^{-11}$	$4.63 \times 10^{-9}$	$9.46 \times 10^{-9}$
	24	$2.47 \times 10^{-2}$	$2.34 \times 10^{-8}$	$9.88 \times 10^{-6}$	$3.25 \times 10^{-7}$
	30	$4.80 \times 10^{-1}$	$3.52 \times 10^{-5}$	$3.47 \times 10^{-2}$	$7.17 \times 10^{-4}$
$q, r = 100$	6	$2.59 \times 10^{-15}$	$2.82 \times 10^{-16}$	$1.20 \times 10^{-11}$	$1.02 \times 10^{-15}$
	12	$4.81 \times 10^{-10}$	$2.94 \times 10^{-14}$	$4.11 \times 10^{-9}$	$1.58 \times 10^{-11}$
	18	$4.33 \times 10^{-5}$	$2.26 \times 10^{-11}$	$1.78 \times 10^{-6}$	$7.83 \times 10^{-9}$
	24	$7.24 \times 10^{-1}$	$2.90 \times 10^{-8}$	$1.37 \times 10^{-2}$	$1.50 \times 10^{-5}$
	30	$3.07 \times 10^{-1}$	$1.45 \times 10^{-5}$	$2.94 \times 10^{-1}$	$4.38 \times 10^{-3}$

Table 8  
Comparison of relative errors in  $x_{1n}$  for Example 17

	$n$	RE_SDA	RE_SDA_m	RE_MSGM	RE_care
$q, r = 1$	6	$1.94 \times 10^{-14}$	$1.11 \times 10^{-15}$	$2.22 \times 10^{-14}$	$5.68 \times 10^{-14}$
	12	$3.99 \times 10^{-10}$	$1.68 \times 10^{-13}$	$4.92 \times 10^{-11}$	$5.61 \times 10^{-11}$
	18	$8.73 \times 10^{-5}$	$6.37 \times 10^{-11}$	$2.35 \times 10^{-8}$	$2.68 \times 10^{-8}$
	24	$3.09 \times 10^{-1}$	$6.39 \times 10^{-8}$	$3.29 \times 10^{-5}$	$6.16 \times 10^{-6}$
	30	$6.49 \times 10^{-1}$	$1.57 \times 10^{-4}$	$1.40 \times 10^{-1}$	$8.37 \times 10^{-3}$
$q, r = 100$	6	$2.13 \times 10^{-15}$	$9.95 \times 10^{-16}$	$3.27 \times 10^{-11}$	$7.67 \times 10^{-15}$
	12	$6.04 \times 10^{-10}$	$1.83 \times 10^{-13}$	$2.30 \times 10^{-8}$	$6.13 \times 10^{-11}$
	18	$5.87 \times 10^{-4}$	$1.16 \times 10^{-10}$	$2.95 \times 10^{-5}$	$2.70 \times 10^{-8}$
	24	$2.02 \times 10^{-1}$	$1.32 \times 10^{-7}$	$2.39 \times 10^{-2}$	$5.20 \times 10^{-5}$
	30	$4.60 \times 10^{-1}$	$5.67 \times 10^{-5}$	$2.98 \times 10^{-1}$	$1.71 \times 10^{-2}$

## 6. Conclusions

Solving CAREs as DAREs, after applying the Cayley transform, has previously been investigated by many. Recent developments and better understanding of doubling algorithms, especially the structure-preserving properties and efficiency of the SDA [24], give this old approach a new lease of life. In addition, we have studied how the parameter  $\gamma$  in the Cayley transform can be chosen optimally. A Fibonacci search for choosing  $\gamma$  was suggested in Section 3, together with the details of other issues involved in the practical implementation of the SDA. We have also developed the SDA\_m which preserves the structure of some doubly symmetric DAREs. Extensive numerical results show that this approach of solving CAREs using the SDA is efficient and competitive, especially for ill-conditioned problems.



## Appendix A

**Lemma A.1.** *If  $\Phi^T = \Phi \geq 0$  and  $\Psi^T = \Psi \geq 0$ , then  $\Phi(I + \Phi\Psi)^{-1} \geq 0$  and  $(I + \Psi\Phi)^{-1} \geq 0$ .*

### Proof

It suffices to prove that  $\Phi(I + \Psi\Phi)^{-1} \geq 0$ . Notice that if  $\Phi$  is positive definite, then the matrix  $\Phi(I + \Psi\Phi)^{-1} = (\Phi^{-1} + \Psi)^{-1}$  is also positive definite. Now, since the matrix  $\Psi^T = \Psi \geq 0$ , we know that  $\Phi + \epsilon I$  is positive definite for  $\epsilon \geq 0$ , and hence we have

$$(\Phi + \epsilon I)[I + \Psi(\Phi + \epsilon I)]^{-1} > 0. \quad (\text{A.1})$$

Let  $\epsilon \rightarrow 0$ , we obtain the desired result. Similarly, it can be shown that  $(I + \Psi\Phi)^{-1}\Psi \geq 0$ .

## Acknowledgment

We would like to thank Professor Ralph Byers and the referee for their valuable comments and suggestions on the manuscript.

## References

- [1] G. Ammar, V. Mehrmann, On Hamiltonian and symplectic Hessenberg forms, *Linear Algebra Appl.* 149 (1991) 55–72.
- [2] B.D.O. Anderson, Second-order convergent algorithms for the steady-state Riccati equation, *Internat. J. Control* 28 (1978) 295–306.
- [3] M. Athans, W. Levine, A. Levis, A system for the optimal and suboptimal position and velocity control for a string of high-speed vehicles, in: *Proc. Fifth Int. Analogue Computation Meetings*, Lausanne, Switzerland, 1967.
- [4] Z. Bai, J. Demmel, On swapping diagonal blocks in real Schur form, *Linear Algebra Appl.* 186 (1993) 73–95.
- [5] Z. Bai, J. Demmel, Using the matrix sign function to compute invariant subspaces, *SIAM J. Matrix Anal. Appl.* 19 (1998) 205–225.
- [6] Z. Bai, Q. Qian, Inverse free parallel method for the numerical solution of algebraic Riccati equations, in: J.G. Lewis (Ed.), *Proc. Fifth SIAM Conf. Appl. Linear Algebra*, Snowbird, UT, June 1994, SIAM, Philadelphia, PA, 1994, pp. 167–171.
- [7] L. Balzer, Accelerated convergence of the matrix sign function, *Internat. J. Control* 21 (1980) 1057–1078.
- [8] A.Y. Barraud, Investigation autour de la fonction signe d'une matrice, application à l'équation de Riccati, *R.A.I.R.O. Automatique* 13 (1979) 335–368.
- [9] A.Y. Barraud, Produit étoile et fonction signe de matrice. application à l'équation de Riccati dans le cas discrete, *R.A.I.R.O. Automatique*, 14 (1980) 55–85.
- [10] M.S. Bazaraa, H.D. Sheraii, C.M. Shetty, *Nonlinear Programming*, John Wiley and Sons, 1993.

- [11] A.N. Beavers, E.D. Denman, Asymptotic solutions to the matrix Riccati equation, *Math. Biosci.* 20 (1974) 339–344.
- [12] A.N. Beavers, E.D. Denman, A computational method for eigenvalues and eigenvectors of a matrix with real eigenvalues, *Numer. Math.* 21 (1974) 389–396.
- [13] A.N. Beavers, E.D. Denman, A new similarity transformation method for eigenvalues and eigenvectors, *Math. Biosci.* 21 (1974) 143–169.
- [14] A.N. Beavers, E.D. Denman, A new solution method for matrix quadratic equations, *Math. Biosci.* 20 (1974) 135–143.
- [15] P. Benner, A.J. Laub, V. Mehrmann, A collection of benchmark examples for the numerical solution of algebraic Riccati equations I: Continuous-time case, Tech. Rep. SPC 95-22, Fakultät für Mathematik, TU Chemnitz-Zwickau, 09107 Chemnitz, FRG, 1995. Available from: <http://www.tu-chemnitz.de/sfb393/spc95pr.html>.
- [16] J.H. Brandts, Matlab code for sorting real Schur forms, *Numer. Linear Algebra Appl.* 9 (2002) 249–261.
- [17] A. Bunse-Gerstner, R. Byers, V. Mehrmann, A chart of numerical methods for structured eigenvalue problems, *SIAM J. Matrix Anal. Appl.* 13 (1992) 419–453.
- [18] A. Bunse-Gerstner, V. Mehrmann, D. Watkins, An SR algorithm for Hamiltonian matrices, based on Gaussian elimination, *Methods Oper. Res.* 58 (1989) 15–26.
- [19] R. Byers, A Hamiltonian QR-algorithm, *SIAM J. Sci. Statist. Comput.* 7 (1986) 212–229.
- [20] R. Byers, Numerical stability and instability in matrix sign function based algorithms, in: C. Byrnes, A. Lindquist (Eds.), *Computational and Combinatorial Methods in System Theory*, North-Holland, 1986, pp. 185–200.
- [21] R. Byers, Solving the algebraic Riccati equation with the matrix sign function, *Linear Algebra Appl.* 85 (1987) 267–279.
- [22] R. Byers, C. He, V. Mehrmann, The matrix sign function method and the computation of invariant subspaces, *SIAM J. Matrix Anal. Appl.* 18 (1997) 615–632.
- [23] E.K.-W. Chu, H.-Y. Fan, W.-W. Lin, A structure-preserving doubling algorithm for continuous-time algebraic Riccati equations, preprint 2002-28, NCTS, National Tsing Hua University, Hsinchu 300, Taiwan, 2003.
- [24] E.K.-W. Chu, H.-Y. Fan, W.-W. Lin, C.-S. Wang, Structure-preserving algorithms for periodic discrete-time algebraic Riccati equations, *Internat. J. Control* 77(8) (2004) 767–788.
- [25] E.K.-W. Chu, H.-Y. Fan, W.-W. Lin, C.-S. Wang, A structure-preserving doubling algorithm for periodic discrete-time algebraic Riccati equations, preprint 2002-18, NCTS, National Tsing Hua University, Hsinchu 300, Taiwan, 2003.
- [26] E.J. Davison, W. Gesing, The systematic design of control systems for the multivariable servomechanism problem, in: M.K. Sain, J.L. Peczkowsky (Eds.), *Alternatives for Linear Multivariable Control*, Nat. Eng. Consortium Inc., Chicago, IL, 1978.
- [27] E. Denman, R. Beavers, The matrix sign function and computations in systems, *Appl. Math. Comput.* 2 (1976) 63–94.
- [28] L. Dieci, Some numerical considerations and Newton's method revisited for solving algebraic Riccati equations, *IEEE Trans. Automat. Control* 36 (1991) 608–616.
- [29] J. Gardiner, A.J. Laub, A generalization of the matrix-sign-function solution to the algebraic Riccati equations, *Internat. J. Control* 44 (1986) 823–832.
- [30] G.H. Golub, C.F. Van Loan, *Matrix Computations*, third ed., The Johns Hopkins University Press, 1996.
- [31] S. Hammarling, Newton's method for solving the algebraic Riccati equation, NPL Rep. DITC 12/82, Nat. Phys. Lab., Teddington, Middlesex TW11 0LW, UK, 1982.
- [32] N.J. Higham, *Accuracy and Stability of Numerical Algorithms*, Society for Industrial and Applied Mathematics, 1996.
- [33] J.L. Howland, The sign matrix and the separation of matrix eigenvalues, *Linear Algebra Appl.* 49 (1983) 221–332.

- [34] P.Hr. Petkov, N.D. Christov, M.M. Konstantinov, On the numerical properties of the Schur approach for solving the matrix Riccati equation, *Systems Control Lett.* 9 (1987) 197–201.
- [35] G.D. Ianculescu, J. Ly, A.J. Laub, P.M. Papadopoulos, Space station freedom solar array  $H_\infty$  control. Talk at 31st IEEE Conf. on Decision and Control, Tucson, AZ, December, 1992.
- [36] M. Kimura, Convergence of the doubling algorithm for the discrete-time algebraic Riccati equation, *Internat. J. Systems Sci.* 19 (1988) 701–711.
- [37] M. Kimura, Doubling algorithm for continuous-time algebraic Riccati equation, *Internat. J. Systems Sci.* 20 (1989) 191–202.
- [38] D. Kleinman, On an iterative technique for Riccati equation computations, *IEEE Trans. Automat. Control* AC-13 (1968) 114–115.
- [39] A.J. Laub, A Schur method for solving algebraic Riccati equations, *IEEE Trans. Automat. Control* 24 (1979) 913–921.
- [40] MathWorks, MATLAB user's guide (for UNIX Workstations), The Math Works, Inc., 1992.
- [41] V. Mehrmann, *The Autonomous Linear Quadratic Control Problem*, Springer-Verlag, 1991.
- [42] V. Mehrmann, A step toward a unified treatment of continuous and discrete time control problems, *Linear Algebra Appl.* 241–243 (1996) 749–779.
- [43] V. Mehrmann, E. Tan, Defect correction methods for the solution of algebraic Riccati equations, *IEEE Trans. Automat. Control* AC-33 (1988) 695–698.
- [44] C.C. Paige, C.F. Van Loan, A Schur decomposition for Hamiltonian matrices, *Linear Algebra Appl.* 41 (1981) 11–32.
- [45] L. Patnaik, N. Viswanadham, I. Sarma, Computer control algorithms for a tubular ammonia reactor, *IEEE Trans. Automat. Control* 25 (1980) 642–651.
- [46] J. Roberts, Linear model reduction and solution of the algebraic Riccati equation by the use of the sign function, *Internat. J. Control* 32 (1980) 667–687.
- [47] N. Sandell, On Newton's method for Riccati equation solution, *IEEE Trans. Automat. Control* AC-19 (1974) 254–255.
- [48] G.W. Stewart, HQR3 and EXCHNG: Fortran subroutines for calculating and ordering the eigenvalues of a real upper Hessenberg matrix, *ACM Trans. Math. Software* 2 (1976) 275–280.
- [49] K. Zhou, J.C. Doyle, K. Glover, *Robust and Optimal Control*, Prentice-Hall, Upper Saddle River, NJ, 1996.