

科技部補助專題研究計畫報告

求解具有高秩常數項的大型代數黎卡迪方程式之數值方法

報告類別：成果報告
計畫類別：個別型計畫
計畫編號：MOST 108-2115-M-003-004-
執行期間：108年08月01日至109年10月31日
執行單位：國立臺灣師範大學數學系（所）

計畫主持人：范洪源

計畫參與人員：碩士班研究生-兼任助理：袁紹宸
碩士班研究生-兼任助理：蔡君彥
碩士班研究生-兼任助理：魏辰亦
碩士班研究生-兼任助理：沈秀宸
碩士班研究生-兼任助理：柯中竣
碩士班研究生-兼任助理：陳世驊
碩士班研究生-兼任助理：蔡宗穎
碩士班研究生-兼任助理：張沛瑀
大專生-兼任助理：許順一
大專生-兼任助理：林以寧

本研究具有政策應用參考價值：否 是，建議提供機關
（勾選「是」者，請列舉建議可提供施政參考之業務主管機關）
本研究具影響公共利益之重大發現：否 是

中華民國 110 年 01 月 23 日

中文摘要：在本專題研究計畫中，我們主要考慮具有高秩常數項的大型尺度代數黎卡迪方程式之數值解。由於此計算解並非數值低秩的情形，因此文獻中各種求取數值低秩解的算法均無法適用於此類問題上，而且目前仍無有效的數值方法解決之。有鑑於此，我們試圖提出修正的doubling演算法求出上述代數黎卡迪方程式的對稱半正定穩定解，並且此算法充分運用了控制輸入矩陣 B 的低秩結構。因為計算解的秩數較高，幾乎接近於滿秩的狀態，以致於計算解的存儲與輸出、殘量與相鄰兩次迭代的相對誤差估計、算法收斂性控制、低秩的迴授增益矩陣輸出等議題的處理手法與過往完全不同。是故，針對這項具有挑戰性的研究課題，吾人已更進一步地探求上述問題的解答。最後，我們亦提供一些有趣且具代表性的數值算例，佐證本計畫所提算法的效率性和可行性。

中文關鍵詞：代數黎卡迪方程式，迴授增益矩陣，高秩常數項，大型尺度問題，LQR最佳控制

英文摘要：In this proposal, we consider the numerical solution of large-scale algebraic Riccati equations with high-rank constant terms. The solutions are not numerically low-rank, so the previously successful methods based on low-rank representations are not directly applicable. We plan to modify the doubling algorithm, making use of the low-rank in the input matrix B . We will also solve the challenging problems in the estimation of residuals and relative errors, convergence control and the output of the modified algorithm. Finally, illustrative numerical examples will be presented for demonstrating the effectiveness and feasibility of the proposed algorithms.

英文關鍵詞：algebraic Riccati equation, feedback gain, high-rank constant term, large-scale problem, LQR optimal control

Large-Scale Algebraic Riccati Equations with High-Rank Constant Terms

Bo Yu* Hung-Yuan Fan[†] Eric King-wah Chu[‡]

Abstract

We consider the numerical solution of large-scale algebraic Riccati equations with high-rank constant terms. The solutions are not numerically low-rank, so the previously successful methods based on low-rank representations are not directly applicable. We modify the doubling algorithm, making use of the low-rank in the input matrix B . We also solve the challenging problems in the estimation of residuals and relative errors, convergence control and the output of the modified algorithm. Illustrative numerical examples are presented.

Keywords. algebraic Riccati equation, feedback gain, high-rank constant term, large-scale problem, LQR optimal control

AMS subject classifications. 15A24, 65F30, 93C05

1 Introduction

One success story of science and technology in the twentieth and twenty first centuries is the ability of mankind to influence or control systems around them. Notable examples of such systems are household appliances, machines, automobiles, aircrafts, computers, the weather, the oceans and human societies. These diverse systems can be approximated or modelled, at least locally, as linear time-invariant (LTI) systems. Algebraic, differential or difference Riccati equations are then solved to construct the optimal controls. Much advances have been made in the solution of these equations in the past few decades. It seems, for a while, that not much more can or needs to be done for algebraic Riccati equations (AREs). This paper proposes the first algorithm for an important and unsolved special case for large-scale AREs with high-rank constant terms.

1.1 Algebraic Riccati Equations

Consider the LTI control system in continuous-time:

$$\dot{x}(t) = Ax(t) + Bu(t), \quad y(t) = Cx(t),$$

*School of Science, Hunan University of Technology, Zhuzhou 412000, China; boyu_hut@126.com

[†]Department of Mathematics, National Taiwan Normal University, Taipei 116, Taiwan; hyfan@math.ntnu.edu.tw (Corresponding author)

[‡]School of Mathematics, Monash University, 9 Rainforest Walk, Melbourne, Victoria 3800, Australia; eric.chu@monash.edu

where $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times m}$ and $C \in \mathbb{R}^{l \times n}$ with $m, l \leq n$. The linear quadratic regulator (LQR) control minimizes the energy or cost functional $J_c(x, u) \equiv \int_0^\infty [x(t)^\top H x(t) + u(t)^\top R u(t)] dt$, with the constant term $H \equiv C^\top T^{-1} C \geq 0$ and $R, T > 0$. Here, a symmetric matrix $P > 0$ (≥ 0) when all its eigenvalues are positive (non-negative). Also, $P > Q$ ($P \geq Q$) if and only if $P - Q > 0$ (≥ 0). The corresponding optimal control $u(t) = F_c x(t)$ and the feedback gain $F_c \equiv -R^{-1} B^\top X$ can then be expressed in terms of the unique positive semidefinite stabilizing solution X of the continuous-time algebraic Riccati equation (CARE) [3, 15, 29, 34, 39], with $G \equiv BR^{-1}B^\top \geq 0$:

$$\mathcal{C}(X) \equiv A^\top X + XA - XGX + H = 0. \quad (1)$$

Analogously, we also consider the LTI control system in discrete-time:

$$x(t+1) = Ax(t) + Bu(t), \quad y(t) = Cx(t).$$

The LQR control minimizes $J_d(x, u) \equiv \sum_{t=0}^\infty [x(t)^\top H x(t) + u(t)^\top R u(t)]$. The corresponding optimal control $u(t) = F_d x(t)$ and the feedback gain $F_d \equiv -(R + B^\top X B)^{-1} B^\top X A$ can be expressed in terms of the unique positive semidefinite stabilizing solution X of the discrete-time algebraic Riccati equation (DARE) [3, 16, 29, 34]:

$$\begin{aligned} \mathcal{D}(X) &\equiv -X + A^\top X(I + GX)^{-1} A + H \\ &= -X + A^\top X A - A^\top X B (R + B^\top X B)^{-1} B^\top X A + H = 0, \end{aligned} \quad (2)$$

by the Sherman-Morrison-Woodbury formula (SMWF) [24]:

$$(M + UDV^\top)^{-1} = M^{-1} - M^{-1}U(D^{-1} + V^\top M^{-1}U)^{-1}V^\top M^{-1}.$$

1.2 Previous Work

The solution of AREs has been an active area of research, due to its importance in optimal control and filtering. Many practitioners in control theory and applied mathematics researched the problem, contributed dozens of methods [15, 16, 29, 34]. Classical methods made use of canonical forms, determinants and polynomial manipulation and state-of-the-art methods solve AREs in a numerically stable manner; see [15, 16] and their references for more details. One favourite approach formulates the AREs as eigenvalue problems [30] and has been implemented in MATLAB [33] (in the commands `care` and `dare`). The other favourite is the Newton-Kleinman method [28]. For modern numerical methods for AREs for systems of moderate dimensions, please consult [5, 6, 8, 9, 15, 16, 30].

For control problems for parabolic PDEs and the balancing based model order reduction of large linear systems, large-scale CAREs, DAREs, Lyapunov and Stein equations have to be solved [2, 3, 5, 9, 11, 12, 13, 25, 26, 27, 35]. Without solving the corresponding AREs, alternative solutions require the invariant or deflating subspaces of the corresponding Hamiltonian matrices or symplectic pencils which are prohibitively expensive to compute.

Benner and his collaborators have done much on large-scale AREs with low-rank structures; see [6, 9, 10, 11, 35] and the references therein. The sub-problems of (inexact) Newton's methods are the Lyapunov and Stein equations which are solved with the ADI iterations. (The initialization of Newton's method and the choice of parameters for the ADI iterations remain challenging.) Consult [1, 26, 27] on the Krylov subspace methods and [18, 19, 20, 21, 22, 23, 38, 37, 40] on the

related gradient based iterations for Lyapunov or Sylvester equations. The structure-preserving doubling algorithm (SDA) [15, 16] has been adapted for large-scale problems [17, 31], making use of the structure in A and the fact that $G \equiv BR^{-1}B^\top$, $H \equiv C^\top T^{-1}C$ are low-rank (i.e. $m, l \ll n$).

1.3 High-Rank H

In many large-scale control problems, the constant term $H \equiv C^\top T^{-1}C$ in the AREs is low-rank with $C \in \mathbb{R}^{l \times n}$ and $l \ll n$. A nice property in this case is that the unique positive definite stabilizing solution X can be approximated numerically by a low-rank matrix $C_a^\top T_a^{-1} C_a$ with $C_a \in \mathbb{R}^{l_a \times n}$ and $l_a \ll n$. Many practitioners then employed the Newton-Galerkin-ADI method [9, 12, 13] or SDA [17, 31] to form C_a and T_a as their computation and storage were efficient.

Recently, Benner *et al* [4, 7, 14] (all unpublished, the only known references for this problem) raised the problem associated with a “high-rank” constant term H . It has been suggested that the (large-scale) CAREs “can be handled using quadrature rules” [7, Page 62], on integrals like

$$K_1 = B^\top X_1 = B^\top \int_0^\infty e^{t(A-BK_0)^\top} (H + K_0^\top K_0) e^{t(A-BK_0)} dt,$$

where $K_0 \equiv B^\top X_0$ is known. A variant of the Newton-Kleinman method is applied when an initial stabilizing X_0 is available.

Unfortunately for a large-scale CARE with a high-rank H , the solution X is no longer numerically low-rank (more details in Sections 2.1 and 2.2). Consequently, the Newton-Galerkin-ADI method or SDA no longer work. Note that even storing or outputting X are nontrivial, so are the computation of the residual or relative errors for any approximation H_k to X and the corresponding convergence control of any iterative procedure for $\{H_k\}$.

1.4 Main Contributions

The efficient solution of large-scale AREs with high-rank constant terms is an important problem which has not been solved. This paper proposes the first algorithm for the problem, which is different to the Newton-Kleinman method in [7], without the difficult initial stabilization step.

Our main contributions are summarized as follows:

- (1) There are some preliminary materials in [4, 14] and some talks (as in [7]) presented on CAREs with high-rank H . Our paper is the first on the numerical solution of the problem.
- (2) We modify the SDA for AREs with high-rank H , revealing the high-rank form of X . Consequently, our algorithm outputs directly the feedback gains $F_d, F_c \in \mathbb{R}^{m \times n}$, which contain only m vectors and are feasible to be represented and stored.
- (3) We show that large-scale AREs for high-rank H can be solved efficiently by the modified doubling algorithm, the SDA_h, which output the feedback gains. We propose algorithms of $O(c_a n)$ computational complexity (with constants involving 2^k , which is in turn $O(1)$ for small values of k). Here $c_a n$ is the flop count required for multiplying a vector v to A or H (or one of their transposes) for DAREs, or the flop count of one solve of the linear system associated with $A - \gamma I$ (or its transpose) for CAREs, and k is the number of iterations required for our (near) quadratic convergent algorithm and is normally small.

- (4) We reveal the new difficulties associated with the computation of the residuals and relative errors, for the convergence control, and propose remedies. These difficulties arise from the high-rank constant term H and are not relevant previously for the solution of AREs with low-rank H in [6, 9, 10, 11, 17, 31, 35].

2 Discrete-Time Algebraic Riccati Equations

Consider the DARE (2) with a *high-rank* constant term $H \geq 0$ but $G = BR^{-1}B^\top \in \mathbb{R}^{n \times n}$ with rank $m \leq n$. For large-scale DAREs, we further assume that $m \ll n$ and A, H are sparse or structured, with at most c_a nonzero elements on each row or column. More generally, we may assume that A, H are “sparse-like”, with multiplication to row or column vectors realizable in $O(c_a n)$ flops. Sparse-like matrices include low-rank updates of banded or sparse matrices.

For the large-scale CARE (1) with $m \ll n$, we need the matrices being efficiently invertible and multiplicable to vectors; see Section 4. CAREs are mathematically equivalent to the DARE (2), through the Cayley transformation (22).

Applying the SDA [16] to the DARE in (2):

$$\begin{cases} G_{k+1} &= G_k + A_k(I + G_k H_k)^{-1} G_k A_k^\top, \\ H_{k+1} &= H_k + A_k^\top H_k (I + G_k H_k)^{-1} A_k, \\ A_{k+1} &= A_k (I + G_k H_k)^{-1} A_k, \end{cases} \quad (3)$$

with $A_0 = A$, $G_0 = G = BR^{-1}B^\top$, $H_0 = H$, we obtain an efficient algorithm for the solution of the DARE with a high rank H . Recall from [16] that $H_k \rightarrow X$, $G_k \rightarrow Y$ (the solution to the dual equation (14) of the DARE (2)) and $A_k \rightarrow 0$, all quadratically ignoring errors. Similar to A_k for large-scale AREs in [17, 31], H_k increases in complexity with respect to k thus cannot be evaluated or stored explicitly, but with its recursion utilized in products with vectors. For large-scale DAREs in [17], both G and H are low-rank and the SMWF can be applied for $(I + G_k H_k)^{-1}$ using the structures in G or H . For our case with a high-rank H , we rely only on the low-rank in G or B . There are complications involving the computation of residuals and convergence control in the SDA, because the high rank of the approximate solution H_k . As elaborated in Section 3, we output instead the approximate feedback gain which is low-rank.

2.1 SDA_h

The details of the SDA_h are described inductively as follows. Consider

$$G_k = B_k R_k B_k^\top, \quad H_k = M_k + C_k T_k C_k^\top, \quad (4)$$

with $B_0 = B$, $R_0 = R^{-1}$, $M_0 = H$, $C_0 = 0$ and $T_0 = 1$. We show (4) holds for all values of k by induction. Note that we used a different convention in [31], where R_k and T_k are written as inverses to emphasize their invertibility. In this paper, R_k , T_k and S_k are positive definite after the truncation and compression in Section 3.1.

First, the application of the SMWF on $I + G_k H_k$ implies that

$$(I_n + G_k H_k)^{-1} = I_n - B_k (I_{m_k} + R_k B_k^\top H_k B_k)^{-1} R_k B_k^\top H_k. \quad (5)$$

From (3) and (5), it is obvious that

$$G_{k+1} = G_k + A_k [I_n - B_k (I_{m_k} + R_k B_k^\top H_k B_k)^{-1} R_k B_k^\top H_k] G_k A_k^\top, \quad (6)$$

with $B_{k+1} = [B_k, A_k B_k]$ and $R_{k+1} = R_k \oplus [R_k - (I_{m_k} + R_k B_k^\top H_k B_k)^{-1} R_k B_k^\top H_k B_k R_k]$. Similarly for A_k , we have

$$A_{k+1} = A_k^2 - D_{k+1}^{(1)} S_{k+1} [D_{k+1}^{(2)}]^\top, \quad (7)$$

with $D_{k+1}^{(1)} = A_k B_k$, $D_{k+1}^{(2)} = A_k^\top H_k B_k$ and $S_{k+1} = (I_{m_k} + R_k B_k^\top H_k B_k)^{-1} R_k$. With $H_k = M_k + C_k T_k C_k^\top$, similar argument implies

$$\begin{aligned} H_{k+1} &= H_k + A_k^\top H_k \left[I_n - B_k (I_{m_k} + R_k B_k^\top H_k B_k)^{-1} R_k B_k^\top H_k \right] A_k \\ &= H_k + A_k^\top H_k A_k - D_{k+1}^{(2)} S_{k+1} [D_{k+1}^{(2)}]^\top, \\ M_{k+1} &= M_k + A_k^\top M_k A_k, \quad C_{k+1} = [C_k, A_k^\top C_k, A_k^\top H_k B_k], \\ T_{k+1} &= T_k \oplus T_k \oplus \left[- (I_{m_k} + R_k B_k^\top H_k B_k)^{-1} R_k \right]. \end{aligned} \quad (8)$$

These imply that A_{k+1} and H_{k+1} inherit the same forms as A_k and H_k in (7) and (8) respectively. From the compression and truncation in Section 3.1, we have B_k being orthogonal.

It is obvious that the full-rank H_k or M_k increase in complexity with increasing k thus cannot be evaluated or stored explicitly. The low-rank matrices $D_k^{(1)}$, $D_k^{(2)}$, S_k , C_k and T_k can be stored and the products of A_k , H_k and M_k with selected vectors can be evaluated recursively.

2.2 Structure of M_k

The structure of M_k has to be explored further because of its importance in our computation. With $M_0 = H$ and LM representing low-rank matrices, we have

$$\begin{aligned} M_1 &= M_0 + A_0^\top M_0 A_0 = H + A^\top H A, \\ M_2 &= M_1 + A_1^\top M_1 A_1 = (H + A^\top H A) + A^\top (I + G H)^{-\top} A^\top (H + A^\top H A) A (I + G H)^{-1} A \\ &= \sum_{j=0}^3 (A^\top)^j H A^j + \text{LM}, \quad \dots \\ M_{k+1} &= M_k + A_k^\top M_k A_k = \sum_{j=0}^{2^k-1} (A^\top)^j H A^j + (A^\top)^{2^k} \left[\sum_{j=0}^{2^k-1} (A^\top)^j H A^j + \text{LM} \right] A^{2^k} + \text{LM} \\ &= \sum_{j=0}^{2^{k+1}-1} (A^\top)^j H A^j + \text{LM}. \end{aligned}$$

The low-rank matrices are easy to handle and can be absorbed into the $C_{k+1} T_{k+1} C_{k+1}^\top$ in H_{k+1} (as in (4)). Consequently, the main problem comes from the first term

$$\check{M}_{k+1} \equiv \sum_{j=0}^{2^{k+1}-1} (A^\top)^j H A^j. \quad (10)$$

For CAREs, A has the form of a shifted inverse and the structure or sparsity in H will not be preserved even for small values of k .

2.3 Stein Equations

Consider the special case of (2) with $G = 0$:

$$\mathcal{S}(X) \equiv -X + A^\top X A + H = 0. \quad (11)$$

The corresponding SDA becomes the Smith method, with the approximate solution H_k equal to M_{k+1} in (10), increasing in complexity. This gives rise to difficulties in convergence control and the evaluation of residuals and successive differences.

In many applications (such as the computation of Gramians in model-order reduction), we ultimately want to estimate the singular values of the solution X of (11). In that situation, the recursion for H_k can be utilized in such an estimation process, in the associated multiplications to some row and column vectors. Similar comments hold for the Lyapunov equations in Section 4.3. Note that this is in the same spirit as in Section 3.3, where the approximations H_k to X are not low-rank but the corresponding approximated feedback gains F_k in (18) and (27) are much simpler (in lower dimensions) thus cheaper to store or output.

3 Computational Issues

3.1 Truncation and Compression

Now we recall an important aspect of the SDA for large-scale DAREs — the growth of B_k . As $H_k = M_k + C_k^\top T_k C_k^\top$ is not computed explicitly and we may control the growth of C_k . We choose not to, because the recursions for A_k and H_k can be applied efficiently enough with only compression on B_k , as seen in (7)–(9).

Obviously, as the SDA converges, increasingly smaller but higher-rank components are added to B_k . Apparent from (6), the growth in the sizes and ranks of these iterates is potentially exponential. If the convergence is slow relative to the growth in B_k , the algorithm will fail.

To reduce the dimensions of B_k , we compress their columns by orthogonalization. As in [17, 31], consider the QR decompositions with column pivoting:

$$B_k = Q_k U_k + \tilde{Q}_k \tilde{U}_k, \quad \|\tilde{U}_k\| \leq \tau_g.$$

From here on, all norms are the 2-norm. Here τ_g is some small tolerance controlling the compression and truncation process, n_k is the number of columns in B_k bounded from above by m_{\max} and its rank satisfies $r_k^{(b)} \equiv \text{rank } B_k \leq n_k \leq m_{\max} \ll n$. Also $Q_k \in \mathbb{R}^{n \times r_k^{(b)}}$ is unitary and $U_k \in \mathbb{R}^{r_k^{(b)} \times n_k}$ is full-rank and upper triangular. We have

$$B_k R_k B_k^\top = Q_k (U_k R_k U_k^\top) Q_k^\top + O(\tau_g), \quad (12)$$

replacing respectively B_k and R_k by the leaner Q_k and $U_k R_k U_k^\top$. We ignore $O(\tau_g)$ terms, control the growth of $r_k^{(b)}$ while sacrificing a hopefully negligible bit of accuracy. We also restrict the widths of B_k , now relabelled $m_k = r_k^{(b)}$ after the compression and truncation in (12), by setting a reasonable upper limit m_{\max} . Similar procedures have been successfully applied in [17, 31].

3.2 Residuals and Differences of Successive Iterates

For the DARE (2) with $H_k = M_k + C_k T_k C_k^\top$, we have the usual relative residual for $H_k \approx X$:

$$\begin{aligned} \mathcal{D}(H_k) &= -H_k + A^\top H_k \left[I_n - B (R + B^\top H_k B)^{-1} B^\top H_k \right] A + H, \\ \hat{r}_k &\equiv \frac{\|\mathcal{D}(H_k)\|}{\|H\| + \|H_k\| + \|A^\top H_k A\| + \|A^\top H_k B (R + B^\top H_k B)^{-1} B^\top H_k A\|}, \end{aligned} \quad (13)$$

analogous to (28), with the help of (2) or the SMWF. This simplifies to $\mathcal{D}(H_k) =$

$$-M_k + A^\top M_k A + H + [-C_k T_k C_k^\top + A^\top C_k T_k C_k^\top A - A^\top H_k B_k (I + R_k B_k^\top H_k B_k)^{-1} R_k B_k^\top H_k A].$$

The norms in (13) are not easy to compute, although can be estimated using the power method [24], which may be attempted after convergence.

Next consider the differences of successive H_k and G_k , respectively:

$$\begin{aligned} \delta H_k &\equiv H_{k+1} - H_k = (M_{k+1} - M_k) + (C_{k+1} T_{k+1} C_{k+1}^\top - C_k T_k C_k^\top) \\ &= A_k^\top M_k A_k + [A_k^\top C_k, A_k^\top H_k B_k] \left[T_k \oplus (I_{m_k} - R_k B_k^\top H_k B_k)^{-1} R_k \right] [A_k^\top C_k, A_k^\top H_k B_k]^\top, \\ \delta G_k &\equiv G_{k+1} - G_k = B_{k+1} R_{k+1} B_{k+1}^\top - B_k R_k B_k^\top = A_k B_k \Delta G_k B_k^\top A_k^\top, \end{aligned}$$

with δG_k being low-rank and $\Delta G_k \equiv R_k - (I_{m_k} + R_k B_k^\top H_k B_k)^{-1} R_k B_k^\top H_k B_k R_k$. Both differences δH_k and δG_k converge to zero quadratically while $A_k \rightarrow 0$ as $k \rightarrow \infty$. Obviously, convergence may be controlled efficiently using δG_k , after $A_k B_k$ has been orthogonalized and δG_k transformed as in Section 3.1, leading to $\|\delta G_k\| = \|\Delta G_k\|$. Note that $\delta G_k \in \mathbb{R}^{n \times n}$ while $\Delta G_k \in \mathbb{R}^{m_k \times m_k}$ with $m_k \ll n$.

3.2.1 Convergence Control and Accuracy of H_k

From the previous Section, it is obvious that the computation of residuals for H_k and the similar quantities is nontrivial and costly. In this section, we describe an alternative of convergence control which bypasses these difficulties. First we quote the following theorem for the SDA:

Theorem 3.1 [32, Theorem 3.1] *Assume that $X, Y > 0$ satisfy the DARE (1) and its dual:*

$$\mathcal{D}_a(Y) \equiv -Y + AY(I + HY)^{-1}A^\top + G = 0, \quad (14)$$

and let $S \equiv (I + GX)^{-1}A$ and $T \equiv (I + HY)^{-1}A^\top$. Then the matrix sequences $\{A_k\}$, $\{G_k\}$ and $\{H_k\}$ generated by the SDA satisfy

- (a) $A_k = (I + G_k X) S^{2^k}$;
- (b) $H \leq H_k \leq H_{k+1} \leq X$ and $X - H_k = (S^\top)^{2^k} (X + X G_k X) S^{2^k} \leq (S^\top)^{2^k} (X + X Y X) S^{2^k}$;
- (c) $G \leq G_k \leq G_{k+1} \leq Y$ and $Y - G_k = (T^\top)^{2^k} (Y + Y H_k Y) T^{2^k} \leq (T^\top)^{2^k} (Y + Y X Y) T^{2^k}$.

Theorem 3.1 implies that the relative accuracy of H_k and G_k are bounded as below:

$$r_k^{(x)} \equiv \frac{\|H_k - X\|}{\|X\|} \leq \|S\|^{2^{k+1}} \|I + XY\|, \quad r_k^{(y)} \equiv \frac{\|G_k - Y\|}{\|Y\|} \leq \|T\|^{2^{k+1}} \|I + XY\|. \quad (15)$$

It is well-known [32] that the spectral radii $\rho(S) = \rho(T)$.

With Theorem 3.1, we can control the convergence of the SDA_h in the following fashion:

- (i) Within an iterative step, compute the (dual) residual $r_k \equiv \|\mathcal{D}_a(G_k)\|$ in $O(n)$ complexity, using the low-rank structure of G_k . Terminate iteration when r_k is small within some tolerance, relative to the data in the dual DARE (as in (28)). This is not that controversial, comparing to the use of $\|A_k\|$ in [15, 16]. The accuracy achieved in H_k is further guaranteed in (ii) and (iii) below.
- (ii) The relative error of H_k is bounded as in (15) with $\|S\| = \|(I + GX)^{-1}A\|$ estimated (with power method [24]). This ignores the effects of any truncation and compression of G_k and H_k , thus underestimates the error.
- (iii) Further indication of the accuracy of H_k can be obtained by applying the power method [24] to estimate δH_k , although only reflecting the relative error of H_{k-1} (thus overestimating the relative error for H_k). This can be costly and should be attempted only after convergence (in the sense of r_k in (i) above).
- (iv) The relative residual \hat{r}_k in (13) is expensive to evaluate. We have applied the power method to estimate the residual after convergence for Example 3 (in Section 5) with $n = 1357$, summarized at the bottom of Table 6. We estimated $\hat{r}_6 = 4.12 \times 10^{-14}$, consistent with the accuracy of other estimates $dH_5 = 8.70 \times 10^{-9}$ and $\phi_6 = 3.72 \times 10^{-27}$. However, this estimation costed more than 4 times the CPU time for the solution of the CARE. The estimation of the relative residuals for higher values of n are even more expensive.
 Note that residuals do not reflect directly on the accuracy of the approximate solutions, without the knowledge of the corresponding condition numbers.
- (v) For CAREs, the theorem and similar techniques can be applied after the Cayley transform.

3.2.2 Estimation of $\|S\|$ and $\|I + XY\|$

After convergence, for $r_k^{(x)}$ in (15), we have

$$\|I + XY\| \approx \|I + H_k G_k\| = \|I + G_k H_k\|, \quad \|S\| \approx \|(I + GH_k)^{-1}A\|.$$

With the help of the SMWF, we have

$$\begin{aligned} \|I + XY\|^2 &\approx \|I + G_k H_k\|^2 = \|I + B_k R_k B_k^\top H_k\|^2 \\ &= \rho(I + B_k R_k B_k^\top H_k + H_k B_k R_k B_k^\top + B_k R_k B_k^\top H_k^2 B_k R_k B_k^\top) = 1 + \rho(\hat{B}_k \hat{R}_k \hat{B}_k^\top), \end{aligned} \quad (16)$$

with

$$\begin{aligned}
\widehat{B}_k &\equiv [B_k, H_k B_k], \quad \widehat{R}_k \equiv \begin{bmatrix} R_k B_k^\top H_k^2 B_k R_k & R_k \\ R_k & 0 \end{bmatrix}, \\
\|S\|^2 &\approx \|(I + GH_k)^{-1}A\|^2 = \|A - B(I + R^{-1}B^\top H_k B)^{-1}R^{-1}B^\top H_k A\|^2 \\
&= \rho \left\{ AA^\top - B(I + R^{-1}B^\top H_k B)^{-1}R^{-1}B^\top H_k AA^\top \right. \\
&\quad \left. - AA^\top H_k B R^{-1}(I + B^\top H_k B R^{-1})^{-1}B^\top \right. \\
&\quad \left. + B(I + R^{-1}B^\top H_k B)^{-1}R^{-1}B^\top H_k AA^\top H_k B R^{-1}(I + B^\top H_k B R^{-1})^{-1}B^\top \right\} \\
&= \rho \left(AA^\top + \check{B}_k \check{R}_k \check{B}_k^\top \right), \tag{17}
\end{aligned}$$

where $\check{R}_k \equiv (I_{m_k} + R^{-1}B^\top H_k B)^{-1}R^{-1}$, $\check{B}_k \equiv [B, AA^\top H_k B]$ and

$$\check{R}_k \equiv \begin{bmatrix} \check{R}_k B^\top H_k AA^\top H_k B \check{R}_k^\top & -\check{R}_k \\ -\check{R}_k^\top & 0 \end{bmatrix}.$$

Consequently, after orthogonalizing \widehat{B}_k in (16) and modifying \widehat{R}_k , we have $\|I + XY\|^2 \approx \|I + G_k H_k\|^2 = 1 + \rho(\widehat{R}_k)$. We also have $\|S\| \approx \|(I + GH_k)^{-1}A\| = [\rho(AA^\top + \check{B}_k \check{R}_k \check{B}_k^\top)]^{1/2}$, involving the spectral radius of a symmetric low-rank update of the sparse matrix AA^\top . The right hand side of (15) can be estimated using the power method [24].

3.3 Feedback Gains

When solving large-scale problems, it is impossible to output the high-rank approximate solution $H_k \in \mathbb{R}^{n \times n}$ explicitly. For large-scale AREs and (non-)linear matrix equations in [17, 31, 36], for instance, the solutions can be represented approximately as numerically low-rank matrices. These can then be stored in an $O(n)$ memory requirement. For our large-scale ARE with a high-rank H , the approximate solution H_k is represented recursively as in (4) and (9) in a nontrivial manner. These difficulties can be bypassed by outputting the estimated feedback gain $F_k \equiv -(R + B^\top H_k B)^{-1}B^\top H_k A \in \mathbb{R}^{m \times n}$, approximating the feedback gain $F_d \equiv -(R + B^\top X B)^{-1}B^\top X A$, for the associated LQR control.

From (4), (7) and (9), we have

$$F_{k+1} = -(R + B^\top H_{k+1} B)^{-1}B^\top H_{k+1} A, \tag{18}$$

with $F_0 = -(R + B^\top H B)^{-1}B^\top H A$. Expanding (18) in terms of F_k does not lead to a more efficient updating formula.

For convergence control, we may utilize

$$\begin{aligned}
\delta F_k &\equiv F_{k+1} - F_k \\
&= -(R + B^\top H_{k+1} B)^{-1}B^\top H_{k+1} A + (R + B^\top H_k B)^{-1}B^\top H_k A \\
&= -(R + B^\top H_{k+1} B)^{-1}B^\top \delta H_k [I_n - B(R + B^\top H_k B)^{-1}B^\top H_k] A.
\end{aligned}$$

3.4 Algorithm and Operation Counts

We summarize the SDA_h below. We only truncate and compress B_{k+1} but with M_{k+1} , T_{k+1} and C_{k+1} in (9) implicitly containing $H_{k+1} = H_k + A_k^\top H_k A_k - D_{k+1}^{(2)} S_{k+1} [D_{k+1}^{(2)}]^\top$. Practically, we only update S_{k+1} , $A_{k+1} B_{k+1}$, $H_{k+1} B_{k+1}$ and $A_{k+1}^\top H_{k+1} V_{k+1}$ with $V_{k+1} = B_{k+1}$ or $A_{k+1} B_{k+1}$.

Algorithm 1 (SDA_h for DAREs with High-Rank H)

Input:	$A, H \in \mathbb{R}^{n \times n}$ and orthogonal $B \in \mathbb{R}^{n \times m}$, $R^{-1} = R^{-\top} \in \mathbb{R}^{m \times m}$, positive tolerances τ_g and ϵ , and m_{\max} ;
Output:	$B_\epsilon \in \mathbb{R}^{n \times m_\epsilon}$, $R_\epsilon \in \mathbb{R}^{m_\epsilon \times m_\epsilon}$, feedback gain $F_\epsilon \in \mathbb{R}^{m \times n}$, dual residual \tilde{r}_ϵ and estimated relative error $r_\epsilon^{(x)}$;
Set	$A_0 = A$, $D_0^{(1)} = D_0^{(2)} = 0 \in \mathbb{R}^{n \times l}$, $S_0 = I_l$, $B_0 = B$, $R_0 = R^{-1}$, $M_0 = H$, $C_0 = 0$ and $T_0 = 1$; $m_0 = m$, $\tilde{r}_0 = 2\epsilon$ and $k = 0$;
Do	until convergence:
If	$\tilde{r}_k \leq \epsilon$, Set $B_\epsilon = B_k$, $R_\epsilon = R_k$ and $\tilde{r}_\epsilon = \tilde{r}_k$;
	Estimate $r_\epsilon^{(x)}$ as in (15); Exit ; End If
Compute	$\hat{B}_{k+1} = [B_k, A_k B_k]$,
	$\hat{R}_{k+1} = R_k \oplus [R_k - (I_{m_k} + R_k B_k^\top H_k B_k)^{-1} R_k B_k^\top H_k B_k R_k]$,
	$D_{k+1}^{(1)} = A_k B_k$, $D_{k+1}^{(2)} = A_k^\top H_k B_k$,
	$S_{k+1} = (I_{m_k} + R_k B_k^\top H_k B_k)^{-1} R_k$;
	with the recursions:
	$A_{k+1} = A_k^2 - D_{k+1}^{(1)} S_{k+1} [D_{k+1}^{(2)}]^\top$,
	$H_{k+1} = H_k + A_k^\top H_k A_k - D_{k+1}^{(2)} S_{k+1} [D_{k+1}^{(2)}]^\top$.
Compute	B_{k+1} by truncating \hat{B}_{k+1} w.r.t. τ_g ;
Set	$R_{k+1} \leftarrow R_{k+1}^\top$, $m_{k+1} \leftarrow \text{rank} B_{k+1}$;
If	$m_{k+1} > m_{\max}$, Set $m_{k+1} \leftarrow m_{\max}$;
	Compute $r_{k+1} \equiv \ D_a(B_{k+1} R_{k+1} B_{k+1}^\top)\ $;
	Compute $q_k \equiv \ G\ + \ G_k\ + \ AG_k(I + HG_k)^{-1} A^\top\ $;
	Compute $\tilde{r}_{k+1} \equiv r_{k+1}/q_k$; End If
	Set $k \leftarrow k + 1$.
End Do	

Recall that the maximum number of nonzero elements in any row or column in A is bounded from above by c_a , for flop counts in the second column in Table 1. In the third column, the number of variables is recorded. Only the dominant operations or memory requirements are included. Most of the work occurs in computing \hat{B}_{k+1} , for which $A_k B_k$ has to be calculated recursively, as A_k and H_k are not available explicitly. We use the notation $N_k \equiv \sum_{j=1}^k m_j$ and a QR decomposition of an $n \times r$ matrix requires $4nr^2$ flops [24, p. 250].

For the details, let s_k and s'_k be the numbers of flops for the products of an arbitrary row or column vectors to A_k and H_k , respectively. After some tedious counting, we deduce that

$$\begin{bmatrix} s_k \\ s'_k \end{bmatrix} = \begin{bmatrix} 2 & 2 \\ 2 + m_{k-1} & 2 \end{bmatrix} \begin{bmatrix} s_{k-1} \\ s'_{k-1} \end{bmatrix} + n \begin{bmatrix} 1 + 5m_{k-1} \\ 0.5m_{k-1}^2 + 4m_{k-1} + 2 \end{bmatrix}.$$

Consequently, for $k \geq 1$, we have $s_k, s'_k \leq \psi_k n$, with $s_0, s'_0 \leq 2c_a m n$ and

$$\begin{aligned} \psi_k &\leq 2(2 + m_{\max})^k c_a m + \frac{1}{4}(m_{\max}^2 + 18m_{\max} + 6) \sum_{j=0}^k (2 + m_{\max})^j \\ &< (2 + m_{\max})^k \left[2c_a m + \frac{k+1}{4}(m_{\max}^2 + 18m_{\max} + 6) \right]. \end{aligned} \quad (19)$$

Table 1: Operation and Memory Counts for the k^{th} Iteration in Algorithm 1 (SDA_h)

Computation	Flops	Memory
\widehat{B}_{k+1}	$m_k s_k$	$N_{k+1} n$
\widehat{R}_{k+1}	$m_k s'_k + 3m_k^2 n$	$O(m_k^2)$
$D_{k+1}^{(1)}, D_{k+1}^{(2)}$	$m_{k+1} [(s_k + s'_k) + 2m_k n]$	—
S_{k+1}	$m_k (s'_k + m_{k+1} n)$	—
B_{k+1}	$4m_k^2 n$	—
R_{k+1}	$O(m_k^3)$	—
\tilde{r}_{k+1}	$8m_k^2 n$	—
Total	$[(m_k + m_{k+1})s_k + (2m_k + m_{k+1})s'_k + 3m_k(5m_k + m_{k+1})n]$	$N_{k+1} n$

Recall that m_k is controlled by the truncation in Section 3.1. In our numerical examples in Section 5, the flop count at the end of Algorithm 1 dominates. This corresponds to the $(2 + m_{\max})^k$ factor in the upper bound of the total flop count. With respect to k , this exponentially growing factor looks menacing. Nonetheless, it behaves as a constant with respect to n , especially with the fast (near) quadratic convergence of the SDA without (with) truncation. This is evident for the examples in Section 5, requiring only 5 to 6 iterations to machine accuracy.

Note that we have not included, in Table 1, the workload for the estimation of ϕ_k (in (30)), $\|\delta H_k\|$ (in (29)) or $\|\widehat{r}_k\|$ (in (13)) by power method. As elaborated later, the associated computation is relatively expensive and should only be conducted selectively after convergence; see Example 3 for the corresponding CPU times.

3.5 Error Analysis

A detailed error analysis for the SDA can be found in [17, Section 3]. Here, we just repeat an abbreviated error analysis for completeness, as from [31, Section 2.4].

Let δG_k , δH_k and δA_k be the errors in G_k , H_k and A_k , respectively, from round-off or the truncation process. Let $\delta_k \equiv \max\{\|\delta G_k\|, \|\delta H_k\|, \|\delta A_k\|\}$ and $\tau \equiv \max\{\tau_g, \tau_h\}$. From the SDA_h in (3) and the truncation process in Section 3.1, we can show that

$$\delta_{k+1} \leq (1 + c_k)(\delta_k + \tau) + O((\delta_k + \tau)^2), \quad (20)$$

with $a_k \equiv \|A_k\|$, $\alpha_k \equiv \|(I + G_k H_k)^{-1}\|$, $\beta_k \equiv \max\{\|G_k\|, \|H_k\|\}$ and $c_k \leq a_k \alpha_k (2\beta_k + 2a_k \alpha_k \beta_k^2 + a_k)$. When the SDA_h converges, we have α_k, β_k converging to constants and a_k to zero as $k \rightarrow \infty$. This indicates that the error δ_k is not magnified through (3). After N iterations, we

have

$$\delta_{N+1} \leq \mathcal{M}_0^{(N)} \delta_0 + \left[\sum_{j=1}^N \mathcal{M}_j^{(N)} \right] \tau + \text{HOTs} \leq \mathcal{M}_0^{(N)} (\delta_0 + N\tau) + \text{HOTs}, \quad (21)$$

with the error magnifying factors $\mathcal{M}_j^{(N)} \equiv \prod_{k=j}^N (1 + c_k)$. Here HOTs denote the higher order terms in δ_k ($k = 0, \dots, N$) and τ . With the SDA_h converges reasonably quickly and $c_k \rightarrow 0$, $\mathcal{M}_0^{(N)}$ will not be too far away from unity. These indicate that the errors in SDA_h behave well and Algorithm 1 is stable in the sense of (20) or (21).

4 Continuous-Time Algebraic Riccati Equations

For CAREs, we assume that linear systems associated with $A - \gamma I$ and its transpose are efficiently solvable in $\tilde{c}_a n$ flops. In other words we assume that \tilde{c}_a is not large. For A being banded, for example, we have $\tilde{c}_a = O(1)$ when Gaussian elimination with partial pivoting is applied.

4.1 Cayley Transform

From [15], the matrices A , G and H in the CARE (1) are first treated with the Cayley transform:

$$A_0 = I + 2\gamma \tilde{A}_\gamma, \quad G_0 = 2\gamma A_\gamma^{-1} G \tilde{A}_\gamma, \quad H_0 = 2\gamma \tilde{A}_\gamma H A_\gamma^{-1}, \quad (22)$$

with $\tilde{A}_\gamma \equiv (A_\gamma + G A_\gamma^{-\top} H)^{-1}$, $A_\gamma \equiv A - \gamma I$ and $\gamma > 0$. The SMWF implies

$$\tilde{A}_\gamma = A_\gamma^{-1} - A_\gamma^{-1} B \cdot (I_m + R^{-1} B^\top A_\gamma^{-\top} H A_\gamma^{-1} B)^{-1} R^{-1} \cdot B^\top A_\gamma^{-\top} H A_\gamma^{-1}. \quad (23)$$

With the above initial A_0 , G_0 and H_0 , the SDA_h still works, again with exactly the same forms and updating formulae for A_k , B_k , $D_k^{(1)}$, $D_k^{(2)}$ and the kernels R_k and S_k . One relevant difference for CAREs is that $A_0 \neq A$ but satisfies, from (22) and (23),

$$A_0 = (I_n + 2\gamma A_\gamma^{-1}) - D_0^{(1)} S_0 \left[D_0^{(2)} \right]^\top, \quad B_0 = A_\gamma^{-1} B. \quad (24)$$

The corresponding rank m perturbed update has the form

$$D_0^{(1)} = B_0, \quad D_0^{(2)} = A_\gamma^{-\top} H B_0, \quad S_0 = 2\gamma (I_m + R^{-1} B_0^\top H B_0)^{-1} R^{-1}. \quad (25)$$

Note that all computations can be realized efficiently, with $A_\gamma^{-1} B$ computed in $O(\tilde{c}_a mn)$ flops.

Similarly, we have $C_0 = 0$, $T_0 = 1$ and

$$R_0 = 2\gamma \left[R^{-1} - R^{-1} B_0^\top H B_0 (I_m + R^{-1} B_0^\top H B_0)^{-1} R^{-1} \right]. \quad (26)$$

4.2 Algorithm and Operation Counts

The initial A_0 , G_0 and H_0 in (22) are different for CAREs, as compared to (3) for DAREs, because of the Cayley transform. These create some differences in the SDA_h for CAREs, which are summarized in this section.

Residuals

For the CARE (1), we have the residual for H_k :

$$\begin{aligned}\check{r}_k &\equiv \mathcal{C}(H_k) = \mathcal{C}(M_k) + [\mathcal{C}(C_k T_k C_k^\top) - H] - M_k B R^{-1} B^\top C_k T_k C_k^\top - C_k T_k C_k^\top B R^{-1} B^\top M_k, \\ \mathcal{C}(M_k) &= A^\top M_k + M_k A - M_k B R^{-1} B^\top M_k + H, \\ \mathcal{C}(C_k T_k C_k^\top) - H &= A^\top C_k T_k C_k^\top + C_k T_k C_k^\top A - C_k T_k C_k^\top B R^{-1} B^\top C_k T_k C_k^\top.\end{aligned}$$

As for DAREs in Section 3.2, the norm of the residual is not easy to compute, although may be estimated using the power method [24] with high additional costs after convergence.

Feedback Gain

For continuous-time, we have the simpler update for the approximate feedback gain:

$$\begin{aligned}F_{k+1} &= -R^{-1} B^\top H_{k+1} = -R^{-1} B^\top (M_{k+1} + C_{k+1} T_{k+1} C_{k+1}^\top) \\ &= -R^{-1} B^\top [(M_k + A_k^\top M_k A_k) + C_{k+1} T_{k+1} C_{k+1}^\top] \\ &= F_k - R^{-1} B^\top (A_k^\top M_k A_k + \check{C}_{k+1} \check{T}_{k+1} \check{C}_{k+1}^\top),\end{aligned}\tag{27}$$

with $\check{C}_{k+1} = [A_k^\top C_k, A_k^\top H_k B_k]$, $\check{T}_{k+1} = T_k \oplus \left[-(I_{m_k} + R_k B_k^\top H_k B_k)^{-1} R_k \right]$ and $F_0 = -R^{-1} B^\top H$. The computation in (18) and (27) can be realized in $O(n)$ complexity.

As for the differences of successive iterates, we have

$$\delta F_k = -R^{-1} B^\top (A_k^\top M_k A_k + \check{C}_{k+1} \check{T}_{k+1} \check{C}_{k+1}^\top).$$

Finally, recall that H_k cannot be explicitly stored or output except in terms of C_k , T_k and the implicit recursion in M_k . These can be used to retrieve the approximate feedback gain F_k .

Operation Counts

The operation counts in Table 1 are still valid for CAREs, with $c_a n$ replaced by $\tilde{c}_a n$, the flop count for one solve of the linear system associated with $A - \gamma I$ or its transpose. There will be an additional start-up cost associated with A_0 , B_0 and R_0 in (22), (24), (25) and (26). This amounts to $4\tilde{c}_a m n$ flops. As for the different costs when A_0 and H_0 are multiplied to row or column vectors, we have to replace c_a in (19) and Table 1 by $2\tilde{c}_a + c_a + 4m_k + 2$. Qualitatively, the operation counts are of similar orders for both DAREs and CAREs, with \tilde{c}_a playing a part for CAREs because of the linear systems associated with A_γ or its transpose.

4.3 Lyapunov Equations

Consider a special case of (1) with $G = 0$, the Lyapunov equation:

$$\mathcal{C}_0(X) \equiv A^\top X + X A + H = 0.$$

As in Section 2.3 for Stein equations, the singular values of X may be required. Multiplications of the approximate solution (in recursion form) to given vectors can then be realized efficiently.

5 Numerical Examples

Here we illustrate the effectiveness of the SDA_h. Only numerical examples with CAREs, which are obviously more difficult to solve because of the Cayley transform involved, are presented here. Algorithm 1 was coded in MATLAB 2010b [33] on a 64-bit PC with 3.4 GHz Intel Core i3 processor and 8G RAM in the Hunan University of Technology. The machine error is reflected by $\text{eps} = 2.22 \times 10^{-16}$ in MATLAB. Recall that $r_k = \|\mathcal{D}_a(G_k)\|$ is the residual for G_k for the dual equation (14). The stopping criterion in terms of the relative residual \tilde{r}_k for G_k is

$$\tilde{r}_k = \frac{r_k}{\|G\| + \|G_k\| + \|AG_k(I + HG_k)^{-1}A^\top\|} \leq \epsilon, \quad (28)$$

for a small tolerance $\epsilon > 0$. We also compute the estimated relative error for H_{k-1} :

$$dH_{k-1} \equiv \frac{\|\delta H_{k-1}\|}{\|H_k\|} = \frac{\|H_{k-1} - H_k\|}{\|H_k\|} \quad (29)$$

and the upper bound of the relative error $r_k^{(x)}$ in (15):

$$\phi_k = \|(I + GH_k)^{-1}A\|^{2^{k+1}} \cdot \|I + G_k H_k\|, \quad (30)$$

estimated by the power method [24] (after convergence). We add the symbol “*” when the power method is terminated prematurely after 150 iterations.

Design of Numerical Examples

Because of the high-rank structure in H (thus H_k and X), the convergence control in (28) is performed making use of the low-rank G_k . For the accuracy of the approximate solution H_k , we have the *estimates* of dH_k and ϕ_k (using the power method) for the relative errors. Examples 1 and 2 are designed with known explicit solutions. Example 2 shows that SDA_h produces approximate solutions at almost the same accuracy as the original SDA in (3) and various estimates of relative errors and residuals are reasonably consistent. Similarly, Examples 2 and 3, modifying a real-life application, illustrate the accuracy and efficiency of the SDA_h.

Example 1. We first tested SDA_h on a simple system with $A = -I_n$, $B = (1, 0, \dots, 0, 1)^\top$, $R = 1$ and $H = 2I_n + BB^\top$ with $n = 10000, 100000, 1000000$. The solution of CARE is $X = I_n$ and the feedback gain $F = -B^\top$. We set the truncation tolerance $\tau_g = 10^{-15}$ for B_k and its maximum width $m_{\max} = 50$. We stopped when $\tilde{r}_k \leq \epsilon = 10^{-15}$. The results are listed in Tables 2–4, where $e_k = \|H_k - I_n\|$ evaluates the true relative accuracy at iteration k , $t_k = \sum_{i=1}^k \delta t_i$ is the sub-total CPU time and δt_i is the CPU-time required for the i^{th} iteration. The relative errors, known exactly, reflect the accuracies of H_k accurately.

Example 2. To assess the performance of SDA_h, we use the real steel profile data from [11], replacing the original matrix A by $A - I_n$. We assign $H = BB^\top - A^\top - A$ for $X = I_n$. In our experiments, we construct explicitly A_{k+1} and H_{k+1} from $D_{k+1}^{(1)}$, S_{k+1} , $D_{k+1}^{(2)}$, A_k and H_k (see also in (7) and (8)). As a comparison, we also computed \hat{H}_{k+1} by the original SDA in (3). Note that only $(n = 1357, m = 7)$ and $(n = 5177, m = 7)$ are considered as the original SDA fails for larger n . The accumulated CPU times of iterations for various dimensions is plotted in Figure 1 and the relative errors $E_h = \|H_k - I\|_2$ and $E_d = \|\hat{H}_k - I\|_2$ are listed Table 5.

k	dH_k	e_k	ϕ_k	r_k	\tilde{r}_k	m_k	δt_k	t_k
1	1.70e-02	2.42e-02	3.13e-02	3.38e-02	3.58e-02	1	0.022	0.022
2	2.48e-02	5.83e-03	1.24e-03	2.19e-03	2.19e-03	2	1.332	1.354
3	4.44e-04	2.28e-05	1.58e-06	8.58e-06	8.58e-06	2	1.331	2.685
4	1.09e-07	1.25e-10	2.49e-12	1.30e-10	1.30e-10	2	3.605	6.290
5	6.52e-15	0.00e-00	6.14e-24	5.75e-17	5.75e-17	2	4.821	11.482

Table 2: Example 1 ($n = 10000$, $\gamma = 1$, $\tau_g = 10^{-15}$, $m_{\max} = 50$)

k	dH_k	e_k	ϕ_k	r_k	\tilde{r}_k	m_k	δt_k	t_k
1	1.70e-02	2.42e-02	3.27e-02	3.38e-02	3.58e-02	1	0.232	0.232
2	2.48e-02	5.83e-03	1.35e-03	2.19e-03	2.19e-03	2	14.702	14.934
3	4.44e-04	2.28e-05	1.88e-06	8.58e-06	8.58e-06	2	14.390	29.324
4	1.09e-07	1.71e-10	3.55e-12	1.30e-10	1.30e-10	2	40.495	69.819
5	6.52e-15	0.00e-00	1.25e-23	5.75e-17	5.75e-17	2	91.299	163.574

Table 3: Example 1 ($n = 100000$, $\gamma = 1$, $\tau_g = 10^{-15}$, $m_{\max} = 50$)

The ranks of the approximate solutions were recorded as m_k and SDA_h stopped when the relative residual $\tilde{r}_k \leq 5.0 \times 10^{-15}$. We see from Figure 1 that SDA_h beats the original doubling algorithm much in CPU time. Table 5 shows that SDA_h attained almost the same accuracy as the SDA and ϕ_5 (for $n = 1357$) or ϕ_6 (for $n = 5177$) were good estimates of the relative errors.

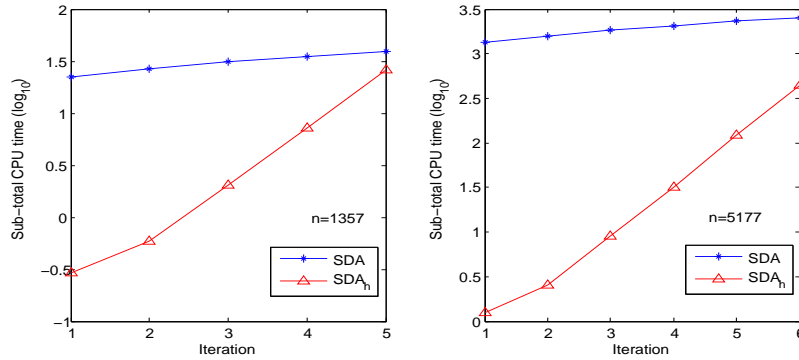


Figure 1: Sub-total CPU time with log for SDA and SDA_h in Example 2.

Example 3. We apply SDA_h on the steel profile example from [11], with A modified to $A - (\|A\|/2)I_n$ and $H = I_n$, for $n = 1357, 5177, 20209, 79841$ and $m = 7$.

The rank m_k was capped at 190 and SDA_h was terminated when $\tilde{r}_k \leq \epsilon = 5.0 \times 10^{-15}$. The accuracy of the approximate solution H_6 achieve high relative residuals of $O(10^{-15})$ efficiently. Adjusting the values of γ (the shift in the Cayley transform from CAREs to DAREs), τ_g and the maximum rank m_{\max} may improve the accuracy of the approximate solution or the speed of convergence further. In our experiments, higher accuracy could also be obtained by setting m_k

k	dH_k	e_k	ϕ_k	r_k	\tilde{r}_k	m_k	δt_k	t_k
1	1.70e-02	2.42e-02	3.31e-02	3.38e-02	3.58e-02	1	2.232	2.232
2	2.48e-02	5.83e-03	1.39e-03	2.19e-03	2.19e-03	2	219.000	221.232
3	4.44e-04	2.28e-05	1.99e-06	8.58e-06	8.58e-06	2	220.473	441.669
4	1.09e-07	2.08e-10	3.96e-12	1.30e-10	1.30e-10	2	616.191	1057.860
5	6.52e-15	0.00e-00	1.57e-23	5.75e-17	5.75e-17	2	1599.382	2695.124

Table 4: Example 1 ($n = 1000000$, $\gamma = 1$, $\tau_g = 10^{-15}$, $m_{\max} = 50$)

$n = 1357$, $\gamma = 4.35$, $\tau_g = 10^{-30}$, $m_{\max} = 190$					
k	E_h	E_d	m_k	r_k	\tilde{r}_k
1	3.65e-02	3.65e-02	14	2.24e-08	3.15e-03
2	1.33e-03	1.33e-03	28	4.05e-10	5.67e-05
3	1.78e-06	1.78e-06	56	2.49e-13	3.49e-08
4	3.19e-12	3.19e-12	112	2.56e-19	3.59e-14
5	1.88e-15	1.88e-15	190	1.06e-20	1.48e-15
$dH_4 = 3.19e-12$			$\phi_5 = 5.15e-22$		
$n = 5177$, $\gamma = 4.35$, $\tau_g = 10^{-30}$, $m_{\max} = 190$					
k	E_h	E_d	m_k	r_k	\tilde{r}_k
1	4.12e-01	4.12e-01	14	7.93e-07	2.31e-02
2	1.70e-01	1.70e-01	28	1.58e-07	4.54e-03
3	2.89e-02	2.89e-02	56	1.06e-08	3.02e-04
4	8.374e-04	8.37e-04	112	9.06e-11	2.58e-06
5	7.01e-07	7.01e-07	190	4.01e-14	1.14e-09
6	4.91e-14	4.91e-14	190	7.36e-20	2.09e-15
$dH_5 = 4.43e-11$			$\phi_6 = 1.19e-23$		

Table 5: Relative error for SDA and SDA_h in Example 2

to a larger positive integer greater than 190, with higher CPU-times and RAM requirements.

The additional CPU-times for the estimation of dH_5 (or \tilde{r}_k) in Tables 3–6 for $n = 1357$, 5177, 20209, 79841 are respectively around 4, 35, 178, 842 (or 5, 38, 182, 854) minutes, with the last two calculations failing to converge within 150 iterations in the power method [24] (but producing reasonably accurate estimates). This illustrates the difficulties mentioned before, in the computation associated with H_k because of its high-rank and recursive structure.

The computation has been controlled by $m_{\max} = 190$ and not the truncation tolerance τ_g (set artificially small). Similar to Example 1, we conclude that dH_k over-estimates the exact relative error e_k while ϕ_k under-estimates. Also, the estimates of dH_5 and \tilde{r}_6 have not converged. Thus H_6 is probably an approximation to an accuracy of $O(10^{-15})$, bounded by the maximum of \tilde{r}_6 and ϕ_6 . The linear progression of t_6 with respect to n indicates an $O(n)$ computational complexity.

6 Conclusions

Large-scale AREs with low-rank constant terms were solved efficiently by the SDA in [17, 31] via the numerically low-rank of the solution X . When the constant terms of AREs are high-rank in some applications [4, 7, 14], calculation of the residual and relative error at each step of the SDA becomes difficult because of the high rank of X . With careful estimation of relative residuals

k	r_k	\tilde{r}_k	m_k	δt_k	t_k
1	1.83e-07	2.57e-02	14	0.184	0.184
2	3.91e-08	5.30e-03	28	0.351	0.535
3	3.59e-09	4.83e-04	56	1.069	1.605
4	4.09e-11	5.49e-06	112	3.795	5.400
5	5.84e-15	7.85e-10	190	14.834	20.234
6	7.63e-21	1.02e-15	190	46.171	66.405
$dH_5 = 8.70e-09$			$\phi_6 = 3.72e-27$ $\hat{r}_6 = 4.76e-15$		

Table 6: Example 3 ($n = 1357, \gamma = 8, \tau_g = 10^{-30}, m_{\max} = 190$)

k	r_k	\tilde{r}_k	m_k	δt_k	t_k
1	2.62e-07	1.87e-02	14	0.866	0.866
2	5.99e-08	4.19e-03	28	1.457	2.323
3	5.10e-09	3.56e-04	56	4.810	7.134
4	4.62e-11	3.22e-06	112	17.894	25.028
5	4.23e-15	2.95e-10	190	67.087	92.116
6	4.02e-20	2.80e-15	190	235.783	327.900
$dH_5 = 8.70e-09$			$\phi_6 = 3.72e-27$ $\hat{r}_6 = 4.12e-14$		

Table 7: Example 3 ($n = 5177, \gamma = 25, \tau_g = 10^{-30}, m_{\max} = 190$)

and errors, we have proposed the SDA_h for large-scale AREs with high-rank constant terms. Efficient control of convergence and economical output of the feedback gain make the SDA_h feasible for CAREs and DAREs. Numerical experiments have validated the effectiveness of the SDA_h.

For future work, we may investigate the possibility of solving large-scale differential and difference equations in X with high-rank constant terms. We propose to output the feedback gains $F_d, F_c \in \mathbb{R}^{m \times n}$, which are much smaller than X when $m \ll n$. We do not know how to tackle the case when $m \approx n$ when n is large. It may well be too greedy, wanting to control a large system with many controls and outputs optimally. After all, we need some simplifying structures to build an efficient algorithm.

Acknowledgements

Part of the work occurred when the first author visited Monash University and the third author visited the School of Mathematical Sciences at Fudan University. The first author was supported partly by the NSF of China (grant 11301170), Natural Science Foundation of Hunan Province (2017JJ2071) and the Excellent Youth Foundation and General Foundation of Hunan Educational Department (17B071, 17C0466).

References

[1] L. AMODEI AND J.-M. BUCHOT, An invariant subspace method for large-scale algebraic Ric-

k	r_k	\tilde{r}_k	m_k	δt_k	t_k
1	5.93e-07	2.17e-02	14	4.723	4.723
2	1.53e-07	5.46e-03	28	8.489	13.213
3	1.76e-08	6.23e-04	56	26.998	40.212
4	3.14e-10	1.11e-05	112	99.008	139.220
5	1.25e-13	4.43e-09	190	392.843	532.064
6	1.98e-19	7.01e-15	190	1490.635	2023.700
$dH_5 = 3.87e-9^*$			$\phi_6 = 2.07e-25$ $\hat{r}_6 = 5.80e-12^*$		

Table 8: Example 3 ($n = 20209$, $\gamma = 108$, $\tau_g = 10^{-30}$, $m_{\max} = 190$)

k	r_k	\tilde{r}_k	m_k	δt_k	t_k
1	1.14e-06	2.10e-02	14	29.045	29.045
2	3.00e-07	5.36e-03	28	50.686	79.728
3	3.31e-08	5.87e-04	56	162.923	242.652
4	5.32e-10	9.43e-06	112	594.442	837.094
5	1.74e-13	3.09e-09	190	2352.497	3189.592
6	2.56e-19	4.45e-15	190	7790.538	10980.130
$dH_5 = 6.36e-9^*$			$\phi_6 = 1.07e-26$ $\hat{r}_6 = 2.95e-10^*$		

Table 9: Example 3 ($n = 79841$, $\gamma = 410$, $\tau_g = 10^{-30}$, $m_{\max} = 190$)

cati equation, *Appl. Numer. Math.*, **60** (2010) 67–1082.

- [2] A. ANTOULAS, *Approximation of Large-Scale Dynamical Systems*, SIAM Publications, Philadelphia, 2005.
- [3] M. ATHAN AND P.L. FALB, *Optimal Control: An Introduction to The Theory and Its Applications*, McGraw-Hill, New York, 1965.
- [4] E. BÄNSCH AND P. BENNER, Stabilization of incompressible flow problems by Riccati-based feedback, In G. Leugering, *et al* (ed.), *Constrained Optimization and Optimal Control for Partial Differential Equations*, International Series of Numerical Mathematics, Vol. **160**, pp. 5–20, Birkhäuser, Basel, 2011; DOI: 10.1007/978-3-0348-0133-1.
- [5] P. BENNER, Solving large-scale control problems, *IEEE Control Syst. Magazine*, **14** (2004) 44–59.
- [6] P. BENNER, Editorial of special issue on “Large-Scale Matrix Equations of Special Type”, *Numer. Lin. Alg. Appl.*, **15** (2008) 747–754.
- [7] P. BENNER, ADI-Based methods for algebraic Lyapunov and Riccati equations, CI-CADA/MIMS Workshop on Numerics for Control and Simulation, Manchester, 2009.
- [8] P. BENNER AND H. FASSBENDER, The symplectic eigenvalue problem, the butterfly form, the SR algorithm, and the Lanczos method, *Lin. Alg. Applic.*, **276** (1998) 19–47.
- [9] P. BENNER AND H. FASSBENDER, On the numerical solution of large-scale sparse discrete-time Riccati equations, *Adv. Comp. Math.*, **35** (2011) 119–147.

- [10] P. BENNER, J.-R. LI AND T. PENZL, Numerical solution of large Lyapunov equations, Riccati equations and linear-quadratic control problems, *Numer. Lin. Alg. Appl.*, **15** (2008) 755–777.
- [11] P. BENNER AND J. SAAK, A semi-discretized heat transfer model for optimal cooling of steel profiles, in *Dimension Reduction of Large-Scale Systems*, ed. P. Benner, V. Mehrmann and D.C. Sorensen, Lecture Notes in Computational Science and Engineering, Springer-Verlag, Berlin, **45** (2005) 353–356.
- [12] P. BENNER AND J. SAAK, A Galerkin-Newton-ADI method for solving large-scale algebraic Riccati equations, *DFG Priority Programme 1253 “Optimization with Partial Differential Equations”*, Preprint **SPP1253-090**, January 2010; available at www.am.uni-erlangen.de/home/spp1253/wiki/images/2/28/Preprint-SPP1253-090.pdf.
- [13] P. BENNER AND J. SAAK, *A Newton-Galerkin-ADI Method for Large-Scale Algebraic Riccati Equations*, Applied Linear Algebra 2010, GAMM Workshop Applied and Numerical Linear Algebra, Novi Sad, May 27, 2010; available at ala2010.pmf.uns.ac.rs/presentations/4g1220pb.pdf.
- [14] P. BENNER AND A. SCHNEIDER, *Balanced Truncation for Descriptor Systems with Many Terminals*, Technical Report, **MPIMD/13-17**, Max Planck Institute Magdeburg, October 2013; available from <http://www.mpi-magdeburg.mpg.de/preprints/>.
- [15] E.K.-W. CHU, H.-Y. FAN AND W.-W. LIN, A structure-preserving doubling algorithm for continuous-time algebraic Riccati equations, *Lin. Alg. Appl.*, **396** (2005) 55–80.
- [16] E.K.-W. CHU, H.-Y. FAN, W.-W. LIN AND C.-S. WANG, A structure-preserving doubling algorithm for periodic discrete-time algebraic Riccati equations, *Int. J. Control*, **77** (2004) 767–788.
- [17] E.K.-W. CHU AND P.C.-Y. WENG, Large-scale discrete-time algebraic Riccati equations — doubling algorithm and error analysis, *J. Comput. Appl. Math.*, **277** (2015) 115–126.
- [18] F. DING AND T.-W. CHEN, Gradient based iterative algorithms for solving a class of matrix equations, *IEEE Trans. Automat. Control*, **50** (2005) 1216–1221.
- [19] F. DING AND T.-W. CHEN, Iterative least squares solutions of coupled Sylvester matrix equations, *Syst. Control Lett.*, **54** (2005) 95–107.
- [20] F. DING AND T.-W. CHEN, On iterative solutions of general coupled matrix equations, *SIAM J. Control Optim.*, **44** (2006) 2269–2284.
- [21] F. DING, P.X. LIU AND J. DING, Iterative solutions of the generalized Sylvester matrix equations by using the hierarchical identification principle, *Appl. Math. Comput.*, **197** (2008) 41–50.
- [22] J. DING, Y. LIU AND F. DING, Iterative solutions to matrix equations of form $A_i X B_i = F_i$, *Comput. Math. Appl.*, **59** (2010) 3500–3507.
- [23] F. DING AND H.-M. ZHANG, Gradient-based iterative algorithm for a class of the coupled matrix equations related to control systems, *IET Control Theory Applic.*, **8** (2014) 1588–1595.

- [24] G.H. GOLUB AND C.F. VAN LOAN, *Matrix Computations*, 3rd ed., Johns Hopkins University Press, Baltimore, 1996.
- [25] M. HEYOUNI AND K. JBILOU, An extended block Arnoldi algorithm for large-scale solutions of the continuous-time algebraic Riccati equations, *Elect. Trans. Numer. Anal.*, **33** (2009) 53–62.
- [26] K. JBILOU, Block Krylov subspace methods for large continuous-time algebraic Riccati equations, *Numer. Algorithms*, **34** (2003) 339–353.
- [27] K. JBILOU, An Arnoldi based algorithm for large algebraic Riccati equations. *Appl. Math. Lett.*, **19** (2006) 437–444.
- [28] D. KLEINMAN, On an iterative technique for Riccati equation computations, *IEEE Trans. Automat. Control*, **13** (1968) 114–115.
- [29] P. LANCASTER AND L. RODMAN, *Algebraic Riccati Equations*, Clarendon Press, Oxford, 1995.
- [30] A.J. LAUB, A Schur method for solving algebraic Riccati equation, *IEEE Trans. Automat. Control*, **AC-24** (1979) 913–921.
- [31] T. LI, E.K.-W. CHU, W.-W. LIN AND P.C.-Y. WENG, Solving large-scale continuous-time algebraic Riccati equations by doubling, *J. Comput. Appl. Math.*, **237** (2013) 373–383.
- [32] W.-W. LIN AND S.-F. XU, Convergence analysis of structure-preserving doubling algorithms for Riccati-type matrix equations, *SIAM J. Matrix Anal. Appl.*, **28** (2006) 26–39.
- [33] MATHWORKS, *MATLAB User's Guide*, 2010.
- [34] V.L. MEHRMANN, *The Autonomous Linear Quadratic Control Problem*, Lecture Notes in Control and Information Sciences, Vol. **163**, Springer Verlag, Berlin, 1991.
- [35] J. SAAK, H. MENA AND P. BENNER, *Matrix Equation Sparse Solvers (MESS): A Matlab Toolbox for the Solution of Sparse Large-Scale Matrix Equations*, Chemnitz University of Technology, 2010.
- [36] P.C.-Y. WENG, H.-Y. FAN AND E.K.-W. CHU, Low-rank approximation to the solution of a nonsymmetric algebraic Riccati equation from transport theory, *Appl. Math. Comput.*, **219** (2012) 729–740.
- [37] L. XIE, Y.-J. LIU AND H.-Z. YANG, Gradient based and least squares based iterative algorithms for matrix equations $AXB + CX^T D = F$, *Appl. Math. Comput.*, **217** (2010) 2191–2199.
- [38] L. XIE, J. DING AND F. DING, Gradient based iterative solutions for general linear matrix equations, *Comput. Math. Appl.*, **58** (2009) 1441–1448.
- [39] C. XU, Y. ZHENG, H. SU AND H. ZENG, Containment for linear multi-agent systems with exogenous disturbances, *Neurocomput.*, **160** (2015) 206–212.
- [40] H. ZHANG AND F. DING, Iterative algorithms for $X + A^T X^{-1} A = I$ by using the hierarchical identification principle. *J. Franklin Inst.*, **353** (2016) 1132–1146.

108年度專題研究計畫成果彙整表

計畫主持人：范洪源		計畫編號：108-2115-M-003-004-			
計畫名稱：求解具有高秩常數項的大型代數黎卡迪方程式之數值方法					
成果項目		量化	單位	質化 (說明：各成果項目請附佐證資料或細項說明，如期刊名稱、年份、卷期、起訖頁數、證號...等)	
國內	學術性論文	期刊論文	0	篇	
		研討會論文	0		
		專書	0	本	
		專書論文	0	章	
		技術報告	0	篇	
		其他	0	篇	
國外	學術性論文	期刊論文	1	篇	B. Yu, H.-Y. Fan, E. K.-W. Chu (2019, Dec). Large-scale algebraic Riccati equations with high-rank constant terms. Journal of Computational and Applied Mathematics, Vol. 361, pp. 130--143. (SCI, 43/260, MATHEMATICS, APPLIED/JCR 2019). MOST 107-2115-M-003-003. 本人為通訊作者.
		研討會論文	0		
		專書	0	本	
		專書論文	0	章	
		技術報告	0	篇	
		其他	0	篇	
參與計畫人力	本國籍	大專生	2	人次	負責收集和整理研究相關資料.
		碩士生	8		負責收集和整理研究相關資料.
		博士生	0		
		博士級研究人員	0		
		專任人員	0		
	非本國籍	大專生	0		
		碩士生	0		
		博士生	0		
		博士級研究人員	0		
		專任人員	0		
其他成果 (無法以量化表達之成果如辦理學術活動、獲得獎項、重要國際合作、研究成果國際影響力及其他協助產業技術發展之具體效益事項等，請以文字敘述填列。)					