

A STRUCTURED DOUBLING ALGORITHM FOR DISCRETE-TIME ALGEBRAIC RICCATI EQUATIONS WITH SINGULAR CONTROL WEIGHTING MATRICES

Chun-Yueh Chiang, Hung-Yuan Fan and Wen-Wei Lin

Abstract. In this paper we propose a structured doubling algorithm for solving discrete-time algebraic Riccati equations without the invertibility of control weighting matrices. In addition, we prove that the convergence of the SDA algorithm is linear with ratio less than or equal $\frac{1}{2}$ when all unimodular eigenvalues of the closed-loop matrix are semi-simple. Numerical examples are shown to illustrate the feasibility and efficiency of the proposed algorithm.

1. INTRODUCTION

The paper concerns with a structured doubling algorithm (SDA) for finding the symmetric almost stabilizing solution X_s of a discrete-time algebraic Riccati equation (DARE) of the form

$$(1.1) \quad \begin{aligned} \mathcal{R}(X) &\equiv -X + A^\top X A + Q \\ &\quad - (C + B^\top X A)^\top (R + B^\top X B)^{-1} (C + B^\top X A) = 0, \end{aligned}$$

where $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times m}$, $C \in \mathbb{R}^{m \times n}$, $Q = Q^\top \in \mathbb{R}^{n \times n}$ and $R = R^\top \in \mathbb{R}^{m \times m}$, respectively. A symmetric solution $X \in \mathbb{R}^{n \times n}$ of (1.1) is called stabilizing (respectively, almost stabilizing) if $R + B^\top X B$ is invertible and all the eigenvalues of the closed-loop matrix $A_F \equiv A + B F$ are in the open (respectively, closed) unit disk, where

$$(1.2) \quad F = -(R + B^\top X B)^{-1} (B^\top X A + C).$$

Moreover, we say that DARE (1.1) has an almost stabilizing solution X with property **(P)** if X is an almost stabilizing solution to DARE (1.1) and all unimodular

Received August 4, 2008, accepted January 8, 2009.

2000 *Mathematics Subject Classification*: 15A24, 47A15.

Key words and phrases: Algebraic Riccati equation, Invariant subspace, Structured doubling algorithm, Singular.

eigenvalues of A_F are semi-simple. The DARE (1.1) arises, e.g., in (a) filtering or stochastic realization problems, and (b) linear quadratic control problems. In case (a), R is the measurement noise covariance and it is not uncommon for this kind of matrix to be singular. For (b), R is the control weighting matrix and in the discrete-time case, such a matrix can occasionally be singular as well. Therefore, we focus on the DARE (1.1) with a singular matrix R throughout the paper.

As is widely known, the DARE (1.1) and its stabilizing solution X_s originated in the discrete-time linear quadratic control problem (case (b) from above) formulated for the discrete-time system

$$(1.3) \quad x_{k+1} = Ax_k + Bu_k, \quad k = 0, 1, \dots, \quad x_0 = \xi.$$

We wish to minimize the cost functional

$$(1.4) \quad J(u) = \sum_{k=0}^{\infty} \begin{bmatrix} x_k^\top & u_k^\top \end{bmatrix} \begin{bmatrix} Q & C^\top \\ C & R \end{bmatrix} \begin{bmatrix} x_k \\ u_k \end{bmatrix},$$

where usually R is positive definite and Q is positive semidefinite (see, e.g., [19] and references therein).

Necessary and sufficient conditions for the existence of the stabilizing solution X_s to the DARE (1.1) are derived in [13] without any assumptions on the invertibility of A or positivity of R or Q . Note that if the DARE (1.1) has a stabilizing solution X_s , then it is unique [13, Proposition 1].

For any $n \times n$ matrices A and B , the matrix pencil $A - \lambda B$ is called regular if $\det(A - \lambda B) \not\equiv 0$. We shall be concerned only with regular pencils throughout the paper. A k -dimensional subspace χ of \mathbb{C}^n is called a deflating subspace for $A - \lambda B$ if there exists matrices $P_1, P_2 \in \mathbb{C}^{n \times k}$ and $Q_1, Q_2 \in \mathbb{C}^{k \times k}$ such that $AP_1 = P_2Q_1$, $BP_1 = P_2Q_2$ and the columns of P_1 span χ . A deflating subspace χ of the pencil $A - \lambda B$ is called stable if the spectrum of $A - \lambda B$ restricted to χ is contained in the open unit disk. On the other hand, a space $\mathcal{V} \in \mathbb{C}^{2n}$ is called isotropic if $x^H \mathcal{J} y = 0$ for any $x, y \in \mathcal{V}$, where the skew-symmetric matrix $\mathcal{J} \equiv \begin{bmatrix} 0 & I \\ -I & 0 \end{bmatrix}$ and I is an $n \times n$ identity matrix. A deflating subspace $\mathcal{V} \subseteq \mathbb{C}^{2n}$ of $A - \lambda B \in \mathbb{C}^{2n \times 2n}$ is called a stable Lagrangian subspace if it is a maximal isotropic subspace and the spectrum of $A - \lambda B$ restricted to \mathcal{V} is contained in the closed unit disk. For solving the symmetric stabilizing solution X_s to the DARE (1.1), one common approach is to compute the stable deflating subspace of the extended symplectic pencil (ESP) $\mathcal{M} - \lambda \mathcal{L}$ associated with the DARE (1.1), where

$$(1.5) \quad \mathcal{M} = \begin{bmatrix} A & 0 & B \\ -Q & I & -C^\top \\ C & 0 & R \end{bmatrix}, \quad \mathcal{L} = \begin{bmatrix} I & 0 & 0 \\ 0 & A^\top & 0 \\ 0 & -B^\top & 0 \end{bmatrix}.$$

This dilated pencil appears naturally when writing the canonical system associated with (1.3) and (1.4) in descriptor form. If DARE (1.1) has a symmetric solution X , then after some elementary block-row operations, we have

$$(1.6) \quad \mathcal{M} - \lambda\mathcal{L} \stackrel{eq.}{\sim} \begin{bmatrix} A + BF - \lambda I & 0 & B \\ 0 & \lambda I - (A + BF)^\top & 0 \\ 0 & \lambda B^\top & R + B^\top X B \end{bmatrix}.$$

According to (1.6), a unimodular number λ is an eigenvalue of $A + BF$ with algebraic multiplicity k if and only if it is an eigenvalue of $\mathcal{M} - \lambda\mathcal{L}$ with algebraic multiplicity $2k$. The following results give some useful properties of the spectrum of ESP.

Lemma 1.1. [8]. *Let λ be a complex number with $|\lambda| = 1$ and X be a solution of (1.1) such that $R + B^\top X B$ is positive definite. If*

$$(1.7) \quad \text{rank}[\lambda I - A, B] = n,$$

then the elementary divisors of $A+BF$ corresponding to λ have degrees k_1, k_2, \dots, k_s ($1 \leq k_1 \leq \dots \leq k_s \leq n$) if and only if the elementary divisors of $\mathcal{M} - \lambda\mathcal{L}$ corresponding to λ have degrees $2k_1, \dots, 2k_s$.

Theorem 1.1. [13]. *Suppose that the ESP (1.5) is regular, then we have:*

1. $\deg \det(\mathcal{M} - \lambda\mathcal{L}) \leq 2m$.
2. *If $\lambda \neq 0$ is a generalized eigenvalue of $\mathcal{M} - \lambda\mathcal{L}$, then $1/\lambda$ is also a generalized eigenvalue of the same multiplicity.*
3. *If $\lambda = 0$ is a generalized eigenvalue of $\mathcal{M} - \lambda\mathcal{L}$ with multiplicity r , then $\lambda = \infty$ is a generalized eigenvalue of multiplicity $m + r$.*

If the stable deflating subspace χ is spanned by columns of a $(2n + m) \times n$ matrix

$$(1.8) \quad V = \begin{bmatrix} V_1 \\ V_2 \\ V_3 \end{bmatrix} \begin{matrix} \}n \\ \}n \\ \}m \end{matrix},$$

and V_1 is invertible, then $X_s = V_2 V_1^{-1}$ is the stabilizing solution of DARE (1.1), see, e.g., [7, 13, 22]. Unfortunately, algorithms based on this property do not take into account the symplectic structure of $\mathcal{M} - \lambda\mathcal{L}$ in (1.5). Non-structure-preserving iterative processes loosen the symplectic structure, and this may cause the algorithm to fail or lose accuracy in adverse circumstances. When $(\mathcal{M}, \mathcal{L})$ has no unimodular eigenvalues and $R > 0$, the quadratically convergent SDA algorithms [6, 11], based

on the viewpoint of the inverse-free iteration [1, 17], have been developed for finding the unique symmetric stabilizing solution X_s of generalized DAREs. Extensive numerical experiments show that the SDA is more efficient and outperforms the other algorithms. Therefore, one of our main purposes is to develop a structured doubling algorithm for computing the symmetric stabilizing solution X_s to DARE (1.1) without the restriction $R > 0$.

Since the DARE (1.1) is a nonlinear matrix equation, it is natural to apply Newton's method (NTM) to obtain an approximate solution. There is an extensive literature on the application of Newton's method for the solution of algebraic Riccati equations, for both the continuous and the discrete case. In the past, an efficient Newton-type method has been proposed by [8] to find the symmetric maximal solution $X_+ \in \mathbb{R}^{n \times n}$ of DARE (1.1). Sufficient conditions for the existence of the maximal solution X_+ to the DARE (1.1) are given in [8, Theorem 1.1].

For $A \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{n \times m}$, the pair (A, B) is said to be d -stabilizable if $\text{rank}[\lambda I - A, B] = n$ for all $\lambda \in \mathbb{C}$ with $|\lambda| \geq 1$. For real symmetric matrices X and Y , we write $X \geq Y$ ($X > Y$) if $X - Y$ is positive semi-definite (definite). A symmetric solution X_+ of (1.1) is called maximal if $X_+ \geq S$ for every symmetric solution S . Maximal and almost stabilizing solutions play important roles in applications, see, e.g., [8, 13, 14, 19] and references given therein. The following theorem tells us that the maximal solution is at least almost stabilizing.

Theorem 1.2. [8]. *Let (A, B) be d -stabilizable pair and assume that there is a symmetric solution \tilde{X} of the inequality $\mathcal{R}(X) \geq 0$ for which $R + B^\top \tilde{X} B > 0$. Then there exists a maximal symmetric solution X_+ of (1.1). Moreover, $R + B^\top X_+ B > 0$ and all the eigenvalues of the closed-loop matrix A_F lie in the closed unit disk.*

As in Theorem 1.4 of [8], Newton's iteration converges quadratically to the symmetric maximal solution X_+ when the same conditions as in Theorem 1.2 are assumed and all eigenvalues of the associated closed-loop matrix A_F are in the open unit disk. In this case, the maximal solution is at least a stabilizing solution of DARE (1.1). Moreover, it has been proven in [8, Theorem 4.3] that Newton's method converges linearly with ratio $\frac{1}{2}$ to the maximal solution X_+ of the DARE (1.1) under the same conditions as in Theorem 1.2 and the ESP in (1.5) satisfies the following assumption.

- (A) All elementary divisors of unimodular eigenvalues of $\mathcal{M} - \lambda\mathcal{L}$ are of degree two.

Assumption (A) is equivalent to the condition that all eigenvalues of the closed-loop matrix $A_F = A - B(R + B^\top X_+ B)^{-1}(C + B^\top X_+ A)$ on the unit circle are semisimple if the first Fréchet derivative of \mathcal{R} in (1.1) at the maximal solution X_+ is not invertible [8] and the conditions of Theorem 1.2 are satisfied. In Section 3,

under the assumption **(A)**, we shall prove that the SDA algorithm converges linearly with ratio less than or equal $\frac{1}{2}$ to an almost stabilizing solution X_s with property **(P)** of the DARE (1.1).

The paper is organized as follows. In Section 2, we propose a structured doubling algorithm for computing the symmetric stabilizing solution or the maximal solution of DARE (1.1) without the invertibility of R . Furthermore, a min-max optimization problem is proposed for selecting a suitable symmetric matrix Y needed to initialize algorithm. The convergence analysis of SDA for solving DAREs (1.1) is shown in Section 3. Some numerical examples with singular matrices R are given in Section 4 to illustrate the efficiency and feasibility of the SDA algorithm. Finally, concluding remarks are given in Section 5.

Throughout this paper, we denote $A^H = \bar{A}^\top$ the conjugate transpose of $A \in \mathbb{C}^{n \times n}$ and $\iota = \sqrt{-1}$. For any positive integer k , I_k and 0_k denote the $k \times k$ identity and zero matrices, respectively. $\|\cdot\|$ denotes any matrix norm, $\rho(A)$ denotes the spectral radius of A .

2. SDA AND NEWTON'S METHOD FOR DARES

In this section we first introduce a structured doubling algorithm for solving the almost stabilizing solution X_s of DARE (1.1) with the control-weighting matrix R being singular. In general, if a symmetric matrix $X \in \mathbb{R}^{n \times n}$ satisfies the DARE (1.1) and all eigenvalues of the closed-loop matrix A_F are in the closed unit disk, then we have

$$(2.1) \quad \mathcal{M} \begin{bmatrix} I \\ X \\ F \end{bmatrix} = \mathcal{L} \begin{bmatrix} I \\ X \\ F \end{bmatrix} \Phi,$$

where the matrix F is as in (1.2) and $\Phi = A_F = A + BF$ with $\rho(\Phi) \leq 1$. Since the control matrix $B \in \mathbb{R}^{n \times m}$ is usually of full column rank in many applications of control system theory, we can select an appropriate matrix $Y = Y^\top \in \mathbb{R}^{n \times n}$ such that $\tilde{R} \equiv R + B^\top Y B$ is invertible. After some elementary block row operators are applied on both sides of (2.1), we obtain

$$(2.2) \quad \begin{bmatrix} (I - G_0 Y)A - B\tilde{R}^{-1}C & 0 & 0 \\ -Q + C^\top \tilde{R}^{-1}(C + B^\top Y A) & I & 0 \\ C + B^\top Y A & 0 & \tilde{R} \end{bmatrix} \begin{bmatrix} I \\ X \\ F \end{bmatrix} \\ = \begin{bmatrix} I - G_0 Y & G_0 & 0 \\ C^\top \tilde{R}^{-1} B^\top Y & A^\top - C^\top \tilde{R}^{-1} B^\top & 0 \\ B^\top Y & -B^\top & 0 \end{bmatrix} \begin{bmatrix} I \\ X \\ F \end{bmatrix} \Phi,$$

where $G_0 = B\tilde{R}^{-1}B^\top$.

Next, post-multiplying the second columns of the matrix pair in (2.2) by Y , and then adding them to the first columns, it follows that

$$(2.3) \quad \begin{bmatrix} (I - G_0Y)A - B\tilde{R}^{-1}C & 0 & 0 \\ \tilde{Q} & I & 0 \\ C + B^\top YA & 0 & \tilde{R} \end{bmatrix} \begin{bmatrix} I \\ X - Y \\ F \end{bmatrix} \\ = \begin{bmatrix} I & G_0 & 0 \\ A^\top Y & A^\top - C^\top \tilde{R}^{-1}B^\top & 0 \\ 0 & -B^\top & 0 \end{bmatrix} \begin{bmatrix} I \\ X - Y \\ F \end{bmatrix} \Phi$$

with $\tilde{Q} = -Q + C^\top \tilde{R}^{-1}(C + B^\top YA) + Y$. Pre-multiplying above matrix pair in (2.3) by the following block elementary matrix

$$\mathcal{E} = \begin{bmatrix} I & 0 & 0 \\ -A^\top Y & I & 0 \\ 0 & 0 & I \end{bmatrix},$$

we thus have

$$(2.4) \quad \begin{bmatrix} A_0 & 0 & 0 \\ -H_0 & I & 0 \\ C + B^\top YA & 0 & \tilde{R} \end{bmatrix} \begin{bmatrix} I \\ X - Y \\ F \end{bmatrix} = \begin{bmatrix} I & G_0 & 0 \\ 0 & A_0^\top & 0 \\ 0 & -B^\top & 0 \end{bmatrix} \begin{bmatrix} I \\ X - Y \\ F \end{bmatrix} \Phi,$$

where

$$(2.5a) \quad A_0 = (I - G_0Y)A - B\tilde{R}^{-1}C,$$

$$(2.5b) \quad G_0 = B\tilde{R}^{-1}B^\top,$$

$$(2.5c) \quad H_0 = A^\top YA_0 - \tilde{Q}.$$

Note that the matrixes G_0 and H_0 are symmetric. Consider the matrix pair $(\mathcal{M}_0, \mathcal{L}_0)$ in standard symplectic form (SSF), where

$$(2.6) \quad \mathcal{M}_0 = \begin{bmatrix} A_0 & 0 \\ -H_0 & I \end{bmatrix}, \quad \mathcal{L}_0 = \begin{bmatrix} I & G_0 \\ 0 & A_0^\top \end{bmatrix}$$

which satisfies $\mathcal{M}_0 \mathcal{J} \mathcal{M}_0^\top = \mathcal{L}_0 \mathcal{J} \mathcal{L}_0^\top$. By Theorem 1.1, Lemma 1.1, (A) and (2.4), it is obvious that the spectrum of $(\mathcal{M}_0, \mathcal{L}_0)$ is the same of $(\mathcal{M}, \mathcal{L})$ except m infinite eigenvalues. The generalized eigenvalues of $(\mathcal{M}_0, \mathcal{L}_0)$ can be arranged as

$$\underbrace{0, \dots, 0}_r, \lambda_{r+1}, \dots, \lambda_\ell, \underbrace{\omega_1, \omega_1, \dots, \omega_{n-\ell}, \omega_{n-\ell}}_{\text{unimodular eigenvalues}}, \overline{\lambda_\ell^{-1}}, \dots, \overline{\lambda_{r+1}^{-1}}, \underbrace{\infty, \dots, \infty}_r,$$

where the eigenvalues λ_i are inside the unit circle except the origin, $i = r+1, \dots, \ell$. From (2.4)–(2.5c), we immediately obtain

$$(2.7) \quad \mathcal{M}_0 \begin{bmatrix} I \\ X - Y \end{bmatrix} = \mathcal{L}_0 \begin{bmatrix} I \\ X - Y \end{bmatrix} \Phi.$$

The DARE associated with the symplectic matrix pair $(\mathcal{M}_0, \mathcal{L}_0)$ in SSF is

$$(2.8) \quad \tilde{X} = A_0^\top \tilde{X} (I + G_0 \tilde{X})^{-1} A_0 + H_0,$$

on which the efficient SDA algorithm [6, 16] can be applied. Note that if \tilde{X} is the symmetric solution to the above DARE (2.8), then $X = \tilde{X} + Y$ is the symmetric solution to the DARE (1.1). As we have mentioned, we can use several elementary block-row and one block column operations to transform the DARE (1.1) to an equivalent DARE (2.8) with the associated symplectic matrix pair $(\mathcal{M}_0, \mathcal{L}_0)$ in SSF. As derived in [6], for given any SSF form (2.6), we construct

$$(2.9) \quad \mathcal{M}_*^{(0)} = \begin{bmatrix} A_0(I + G_0 H_0)^{-1} & 0 \\ -A_0^\top(I + H_0 G_0)^{-1} H_0 & I \end{bmatrix}, \quad \mathcal{L}_*^{(0)} = \begin{bmatrix} I & A_0 G_0 (I + H_0 G_0)^{-1} \\ 0 & A_0^\top (I + H_0 G_0)^{-1} \end{bmatrix}$$

and consequently deduce that

$$(2.10) \quad \mathcal{M}_*^{(0)} \mathcal{L}_0 = \mathcal{L}_*^{(0)} \mathcal{M}_0.$$

We now compute $\mathcal{L}_*^{(0)} \mathcal{L}_0$ and $\mathcal{M}_*^{(0)} \mathcal{M}_0$, and apply the Sherman-Morrison-Woodbury formula to produce

$$(2.11) \quad \mathcal{M}_1 \equiv \mathcal{M}_*^{(0)} \mathcal{M}_0 = \begin{bmatrix} A_1 & 0 \\ -H_1 & I \end{bmatrix}, \quad \mathcal{L}_1 \equiv \mathcal{L}_*^{(0)} \mathcal{L}_0 = \begin{bmatrix} I & G_1 \\ 0 & A_1^\top \end{bmatrix},$$

where

$$(2.12a) \quad A_1 = A_0(I + G_0 H_0)^{-1} A_0,$$

$$(2.12b) \quad G_1 = G_0 + A_0 G_0 (I + H_0 G_0)^{-1} A_0^\top,$$

$$(2.12c) \quad H_1 = H_0 + A_0^\top (I + H_0 G_0)^{-1} H_0 A_0.$$

Equations in (2.11) show that the matrix pair $(\mathcal{M}_1, \mathcal{L}_1)$ is again in SSF form. From (2.10)–(2.11), the pair $(\mathcal{M}_1, \mathcal{L}_1)$ satisfies the doubling property: if $\mathcal{M}_0 x = \lambda \mathcal{L}_0 x$, then $\mathcal{M}_1 x = \lambda \mathcal{M}_*^{(0)} \mathcal{L}_0 x = \lambda \mathcal{L}_*^{(0)} \mathcal{M}_0 x = \lambda^2 \mathcal{L}_*^{(0)} \mathcal{L}_0 x = \lambda^2 \mathcal{L}_1 x$. Equations (2.12a)–(2.12c) form the basis of the SDA [6, 16], which can be modified for the DAREs (1.1) as follows:

Algorithm 2.1 (SDA for DAREs).

```

Input:  $A, B, C, Q, R; \tau$  (a small tolerance);  $k=0, err=1$ ;
Output: a symmetric stabilizing solution  $X$  to DARE (1.1).
Select a symmetric matrix  $Y$  such that  $\tilde{R} \equiv R + B^\top Y B$  is invertible;
Put  $A_0 := (I - GY)A - B\tilde{R}^{-1}C$ ,
 $G_0 := B\tilde{R}^{-1}B^\top$ ,
 $H_0 := Q - Y - C^\top \tilde{R}^{-1}B^\top Y A - A^\top Y B\tilde{R}^{-1}C$ 
 $- C^\top \tilde{R}^{-1}C + A^\top Y(I - GY)A$ ;
While  $err > \tau$ ,
Put  $A_{k+1} := A_k(I + G_k H_k)^{-1} A_k^\top$ ,
 $G_{k+1} := G_k + A_k G_k (I + H_k G_k)^{-1} A_k^\top$ ,
 $H_{k+1} := H_k + A_k^\top (I + H_k G_k)^{-1} H_k A_k$ ,
 $err := \frac{\|H_{k+1} - H_k\|}{\max\{1, \|H_k\|\}}$ ;
If  $I + G_k H_k$  is ill-conditioned, then break down,
Else set  $k := k + 1$ ;
End if
End
 $X := H_k + Y$ .

```

The Newton's Method in [8, 14] is developed to solve the DARE (1.1) by solving a discrete-time Lyapunov equation (or Stein) at each iteration. The convergence of Newton's method is shown to be either quadratic or linear with the common ratio $\frac{1}{2}$. Specifically, the Newton's method can be stated as follows. Here we use the Matlab command `dlyap` to solve the Stein equation [18].

Algorithm 2.2 (NTM for DAREs).

```

Input:  $A, B, C, Q, R; \tau$  (a small tolerance);  $k=0, err=1$ ;
Output: a symmetric stabilizing solution  $X$  to DARE.
Choose  $L_0$  such that  $A_0 \equiv A - B L_0$  is  $d$ -stable;
Solve  $X_0 := \text{dlyap}(A_0^\top, Q + L_0^\top R L_0 - C^\top L_0 - L_0^\top C)$ ;
While  $err > \tau$ ,
Put  $L_{k+1} := (R + B^\top X_k B)^{-1} (C + B^\top X_k A)$  and  $A_{k+1} := A - B L_k$ ;
Solve  $X_{k+1} := \text{dlyap}(A_{k+1}^\top, Q + L_{k+1}^\top R L_{k+1} - C^\top L_{k+1} - L_{k+1}^\top C)$ ;
Put  $err := \frac{\|X_{k+1} - X_k\|}{\max\{1, \|X_k\|\}}$ ;
Set  $k := k + 1$ ;
End
 $X := X_k$ .

```

2.1. Selection of Y

For simplicity, we choose $Y = \gamma I \in \mathbb{R}^{n \times n}$ such that $R_\gamma \equiv R + \gamma B^\top B$ is invertible for $\gamma > 0$. We first derive the forward error bounds of matrices A_0, G_0

and H_0 given in (2.5a)–(2.5c), 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 $fl(\cdot)$ to denote computed floating point values. The quantity 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 $A \preceq B$ if $a_{ij} \leq b_{ij}$ for all i, j . The 1- and ∞ - matrix norms are denoted by $\|\cdot\|_1$ and $\|\cdot\|_\infty$, respectively.

Since $Y = \gamma I$, it follows from (2.5a)–(2.5c) that

$$(2.13) \quad A_0 = A - \gamma BR_\gamma^{-1}B^\top A - BR_\gamma^{-1}C,$$

$$(2.14) \quad G_0 = BR_\gamma^{-1}B^\top,$$

$$(2.15) \quad \begin{aligned} H_0 &= Q - C^\top R_\gamma^{-1}C - \gamma I - \gamma C^\top R_\gamma^{-1}B^\top A - \gamma A^\top BR_\gamma^{-1}C \\ &\quad + \gamma A^\top A - \gamma^2 A^\top BR_\gamma^{-1}B^\top A. \end{aligned}$$

Since $R_\gamma \in \mathbb{R}^{m \times m}$ is symmetric indefinite, the matrix $W_B \equiv R_\gamma^{-1}B^\top$ can be computed by block LDL[⊤] factorization with any pivoting strategy, for instance, the Bunch-Kaufman partial pivoting strategy, see e.g., [10, Chapter 11]. Suppose this algorithm yields the computed factorization $PR_\gamma P^\top \approx \tilde{L}\tilde{D}\tilde{L}^\top$, where P is a permutation matrix and \tilde{D} has diagonal blocks of dimension 1 or 2. From Theorem 11.3 of [10], we obtain

$$(2.16) \quad \begin{aligned} fl(W_B) &= W_B + E_1, \\ |E_1| &\leq p_1(m)u \left[|R_\gamma^{-1}|(|R_\gamma| + P^\top|\tilde{L}||\tilde{D}||\tilde{L}^\top|P)|fl(W_B)| \right] + O(u^2), \end{aligned}$$

where $p_1(m)$ is a linear polynomial. Since it can be shown that the matrix $|\tilde{L}||\tilde{D}||\tilde{L}^\top|$ satisfies the bound [9]

$$(2.17) \quad \|\tilde{L}||\tilde{D}||\tilde{L}^\top\|_M \leq 36m\rho_m\|R_\gamma\|_M,$$

where $\|R_\gamma\|_M \equiv \max_{i,j} |(R_\gamma)_{ij}|$ and ρ_m is the growth factor.

Therefore, it follows from (2.16) and (2.17) that the forward error E_1 satisfies

$$(2.18) \quad \begin{aligned} \|E_1\|_\infty &\leq p_1(m)u \left[\|R_\gamma^{-1}\|_\infty\|R_\gamma\|_\infty\|fl(W_B)\|_\infty \right. \\ &\quad \left. + \|R_\gamma^{-1}\|_\infty\|\tilde{L}||\tilde{D}||\tilde{L}^\top\|_\infty\|fl(W_B)\|_\infty \right] + O(u^2) \\ &\leq p_1(m)u \left[\|R_\gamma^{-1}\|_\infty\|R_\gamma\|_\infty\|fl(W_B)\|_\infty \right. \\ &\quad \left. + m\|R_\gamma^{-1}\|_\infty\|\tilde{L}||\tilde{D}||\tilde{L}^\top\|_M\|fl(W_B)\|_\infty \right] + O(u^2) \end{aligned}$$

$$\begin{aligned} &\leq p_1(m)(1 + 36m^2\rho_m)\rho_mu(\|R_\gamma^{-1}\|_\infty\|R_\gamma\|_\infty\|fl(W_B)\|_\infty) + O(u^2) \\ &\leq p(m)\rho_mu(\|R_\gamma^{-1}\|_\infty\|R_\gamma\|_\infty\|fl(W_B)\|_\infty) + O(u^2), \end{aligned}$$

where $p(m)$ is a cubic polynomial. Similarly, the forward error bound in evaluating $W_C \equiv R_\gamma^{-1}C$ is given by

$$(2.19) \quad \begin{aligned} fl(W_C) &= W_C + E_2, \\ \|E_2\|_\infty &\leq p(m)\rho_mu(\|R_\gamma^{-1}\|_\infty\|R_\gamma\|_\infty\|fl(W_C)\|_\infty) + O(u^2). \end{aligned}$$

Furthermore, it can be derived from (2.18) and (2.19) that

$$(2.20) \quad \begin{aligned} fl(\gamma BR_\gamma^{-1}B^\top A) &= \gamma BR_\gamma^{-1}B^\top A + E_3, \\ \|E_3\|_\infty &\leq p(m)\rho_mu(\gamma\|B\|_\infty\|R_\gamma^{-1}\|_\infty\|R_\gamma\|_\infty\|fl(W_B)\|_\infty\|A\|_\infty) \\ &\quad + O(u^2), \end{aligned}$$

and

$$(2.21) \quad \begin{aligned} fl(BR_\gamma^{-1}C) &= BR_\gamma^{-1}C + E_4, \\ \|E_4\|_\infty &\leq p(m)\rho_mu(\|B\|_\infty\|R_\gamma^{-1}\|_\infty\|R_\gamma\|_\infty\|fl(W_C)\|_\infty) \\ &\quad + mu(\|B\|_\infty\|R_\gamma^{-1}\|_\infty\|C\|_\infty) + O(u^2). \end{aligned}$$

Therefore, from (2.18), (2.20) and (2.21), we can deduce that the forward error bounds in evaluating A_0 and G_0 in (2.13)-(2.14) are

$$(2.22) \quad \begin{aligned} fl(A_0) &= A_0 + E_5, \\ \|E_5\|_\infty &\leq p(m)\rho_mu [(\gamma\|A\|_\infty + 1)\|B\|_\infty\kappa_\infty(R_\gamma)\|fl(W_B)\|_\infty] \\ &\quad + (2n + 3)u(\gamma\|B\|_\infty\|R_\gamma^{-1}\|_\infty\|B\|_1\|A\|_\infty) \\ &\quad + (m + 1)u(\|B\|_\infty\|R_\gamma^{-1}\|_\infty\|C\|_\infty) \\ &\quad + 2u\|A\|_\infty + O(u^2), \end{aligned}$$

and

$$(2.23) \quad \begin{aligned} fl(G_0) &= G_0 + E_6, \\ \|E_6\|_\infty &\leq p(m)\rho_mu(\|B\|_\infty\kappa_\infty(R_\gamma)\|fl(W_B)\|_\infty) \\ &\quad + mu(\|B\|_\infty\|R_\gamma^{-1}\|_\infty\|B\|_1) + O(u^2), \end{aligned}$$

where $\kappa_\infty(R_\gamma) = \|R_\gamma^{-1}\|_\infty\|R_\gamma\|_\infty$ is the condition number of R_γ .

Finally, from (2.18) and (2.19), the forward error bound in evaluating the matrix H_0 in (2.15) is given by

$$\begin{aligned}
 fl(H_0) &= H_0 + E_7, \\
 \|E_7\|_\infty &\leq p(m)\rho_m u(\|C\|_1 \kappa_\infty(R_\gamma) \|fl(W_C)\|_\infty \\
 &\quad + \gamma \|C\|_1 \kappa_\infty(R_\gamma) \|fl(W_B)\|_\infty \|A\|_\infty \\
 &\quad + \gamma \|A\|_1 \|B\|_\infty \kappa_\infty(R_\gamma) \|fl(W_C)\|_\infty \\
 &\quad + \gamma^2 \|A\|_1 \|B\|_\infty \kappa_\infty(R_\gamma) \|fl(W_B)\|_\infty \|A\|_\infty) \\
 (2.24) \quad &+ (m+6)u(\|C\|_1 \|R_\gamma^{-1}\|_\infty \|C\|_\infty) \\
 &\quad + 6u(\|Q\|_\infty + \gamma) + (n+3)u(\gamma \|A\|_1 \|A\|_\infty) \\
 &\quad + (n+m+5)u(\gamma \|C\|_1 \|R_\gamma^{-1}\|_\infty \|B\|_1 \|A\|_\infty) \\
 &\quad + (n+m+4)u(\gamma \|A\|_1 \|B\|_\infty \|R_\gamma^{-1}\|_\infty \|C\|_\infty) \\
 &\quad + (3n+2)u(\gamma^2 \|A\|_1 \|B\|_\infty \|R_\gamma^{-1}\|_\infty \|B\|_1 \|A\|_\infty) + O(u^2).
 \end{aligned}$$

In order to control the forward error bounds of A_0 , G_0 and H_0 , and the conditioning of $I + G_0 H_0$, we consider the following min-max optimization problem, to determine an optimal value $\hat{\gamma} > 0$:

$$(2.25) \quad \min_{\gamma > 0} F(\gamma), \quad F(\gamma) \equiv \max\{f_i(\gamma), i = 1, 2, 3\},$$

where the functions $f_1(\gamma) = \kappa_\infty(R_\gamma)$, $f_2(\gamma) = \gamma^2 \kappa_\infty(R_\gamma)$ and $f_3(\gamma) = \text{cond}(I + G_0 H_0)$, respectively.

Since the condition number $\kappa_\infty(R_\gamma)$ is bounded 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 $F(\gamma)$ in Figure 2.1.

We can apply the Fibonacci search method to compute an approximate value of $\hat{\gamma}$, see, e.g., [3, p. 272]. In order to save computational costs, our experience indicates that three to five iterations of Fibonacci search are adequate to obtain a suboptimal yet acceptable approximation to $\hat{\gamma}$.

3. CONVERGENCE OF SDA

In [13], necessary and sufficient conditions are given for the existence of the stabilizing solution $X_s = X_s^\top \in \mathbb{R}^{n \times n}$ to DARE (1.1).

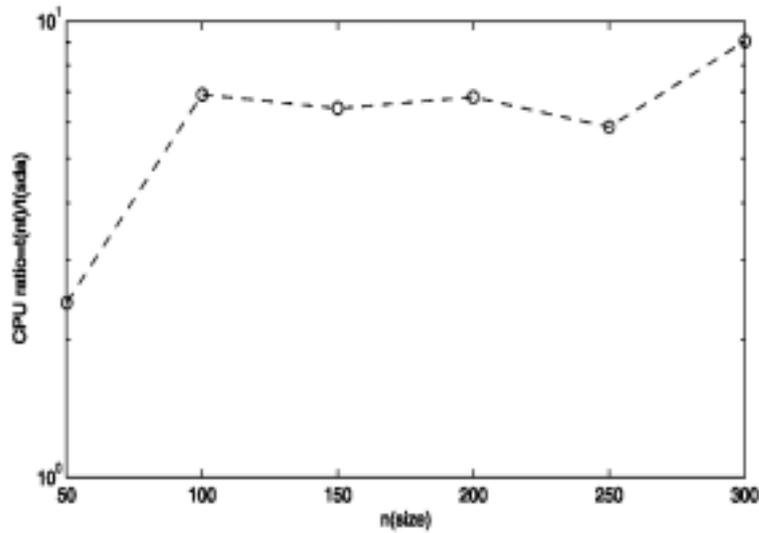


Fig. 2.1. The graph of $F(\gamma)$.

Definition 3.1. [13]. A regular ESP is called *disconjugate* if it has no generalized eigenvalues on the unit circle and V_1 is invertible in (1.8).

It has been proven in [13] that the DARE (1.1) has a unique stabilizing solution X_s if and only if the ESP is disconjugate. In this section we shall first characterize the quadratic convergence of SDA under the same conditions. It is easily seen from (2.7) that the standard symplectic pencil $\mathcal{M}_0 - \lambda\mathcal{L}_0$ in (2.6) is also disconjugate if the ESP is disconjugate. Therefore, we deduce that the simplified DARE (2.8) has a unique stabilizing solution \tilde{X}_s . For simplicity, we only consider complex matrices in the following convergence theorems. The proof for real symplectic pencils can be modified slightly from the complex cases. Suppose that there exist nonsingular U, W such that

$$(3.1) \quad U\mathcal{M}_0W = \begin{bmatrix} J_s & 0 \\ 0 & I \end{bmatrix}, \quad U\mathcal{L}_0W = \begin{bmatrix} I & 0 \\ 0 & J_s \end{bmatrix},$$

where J_s is the stable Jordan block of size n , i.e., $\rho(J_s) < 1$. If we denote

$$(3.2) \quad W = \begin{bmatrix} W_1 & W_3 \\ W_2 & W_4 \end{bmatrix},$$

where $W_i \in \mathbb{C}^{n \times n}$ for $i = 1, 2, 3, 4$, the quadratic convergence of the SDA algorithm has been proved in [6, Theorem 1].

Theorem 3.1. [6]. *Suppose that the ESP is disconjugate. If W_1 and W_4 in (3.2) are invertible, then the sequences $\{A_k, H_k, G_k\}$ computed by the SDA algorithm satisfy:*

- (i) $\|A_k\| = O(\|J_s^{2k}\|) \rightarrow 0$ as $k \rightarrow \infty$.
- (ii) $H_k \rightarrow \tilde{X}_s$, where \tilde{X}_s is a stabilizing solution of DARE (2.8)

$$\tilde{X} = A_0^\top \tilde{X} (I + G_0 \tilde{X})^{-1} A_0 + H_0.$$

- (iii) $G_k \rightarrow \tilde{X}_d$, where \tilde{X}_d solves the dual DARE

$$\tilde{Y} = A_0 \tilde{Y} (I + H_0 \tilde{Y})^{-1} A_0^\top + G_0.$$

Moreover, the convergence rate in (i)–(iii) above is $O(|\lambda_n|^{2k})$, 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 $\mathcal{M}_0 - \lambda \mathcal{L}_0$ (including 0 and ∞).

Remark 3.1. *Assuming the conditions in Theorem 3.1 and \tilde{X}_s is the unique stabilizing solution of DARE (2.8), it follows that the symmetric matrix $X_s = \tilde{X}_s + Y$ must be the unique maximal, stabilizing solution of DARE (1.1).*

On the other hand, when the ESP satisfies the condition (A), we shall prove the linear convergence of SDA with ratio less than or equal to $\frac{1}{2}$. Denote the Jordan block of size p corresponding to a unimodular eigenvalue $\omega \equiv e^{i\theta}$ by

$$(3.3) \quad J_{\omega,p} = \begin{bmatrix} \omega & 1 & 0 & \cdots & 0 \\ 0 & \omega & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ \vdots & & & \ddots & 1 \\ 0 & \cdots & \cdots & 0 & \omega \end{bmatrix}_{p \times p}.$$

In particular, for the unimodular eigenvalues $\omega_j = e^{i\theta_j}$ of the matrix pair $(\mathcal{M}_0, \mathcal{L}_0)$ with $p = 1$, we have

$$(3.4) \quad J_{\omega_j,1} = [\omega_j]$$

for $j = 1, \dots, n - \ell$. From symplectic Kronecker Theorem for $(\mathcal{M}, \mathcal{L})$ (see [15]), there exist a symplectic matrix \mathcal{Z} (i.e., $\mathcal{Z}^\top \mathcal{J} \mathcal{Z} = \mathcal{J}$) and a nonsingular \mathcal{Q} such that

$$(3.5a) \quad \mathcal{Q} \mathcal{M}_0 \mathcal{Z} = \begin{bmatrix} J_s \oplus J_1 & 0_\ell \oplus \Gamma_1 \\ 0 & I_\ell \oplus J_1^{-H} \end{bmatrix} \equiv \mathcal{J}_\mathcal{M},$$

$$(3.5b) \quad \mathcal{Q} \mathcal{L}_0 \mathcal{Z} = \begin{bmatrix} I_\ell \oplus I_\mu & 0 \\ 0 & J_s^H \oplus I_\mu \end{bmatrix} \equiv \mathcal{J}_\mathcal{L},$$

where J_s is the stable Jordan block of size ℓ ,

$$(3.6) \quad J_1 = J_{\omega_1,1} \oplus \cdots \oplus J_{\omega_\mu,1},$$

$$(3.7) \quad \Gamma_1 = I_\mu,$$

with $\mu = n - \ell$ and \oplus denoting the direct sum of matrices. Note that $J_1^{-H} = J_1$, and the matrices \mathcal{J}_M and \mathcal{J}_L in (3.5) commute with each other. Therefore, from (3.5), one can derive

$$(3.8) \quad \mathcal{M}_0 \mathcal{Z} \mathcal{J}_L = \mathcal{Q}^{-1} \mathcal{J}_L \mathcal{J}_M = \mathcal{L}_0 \mathcal{Z} \mathcal{J}_M.$$

From (3.5) it follows that $\text{span}\{\mathcal{Z}(:, 1 : n)\}$ forms the unique stable Lagrangian deflating subspace of $(\mathcal{M}, \mathcal{L})$ corresponding to $J_s \oplus J_1$. Let $\{(\mathcal{M}_k, \mathcal{L}_k)\}_{k=0}^\infty$ be the sequence of symplectic pairs in SSF with

$$(3.9) \quad \mathcal{M}_k = \begin{bmatrix} A_k & 0 \\ -H_k & I \end{bmatrix}, \quad \mathcal{L}_k = \begin{bmatrix} I & G_k \\ 0 & A_k^\top \end{bmatrix}$$

generated by Algorithm 2.1. It follows from (3.8) as well as (2.10)–(2.11) that

$$(3.10) \quad \begin{aligned} \mathcal{M}_1 \mathcal{Z} \mathcal{J}_L^2 &= \mathcal{M}_*^{(0)} \mathcal{M}_0 \mathcal{Z} \mathcal{J}_L^2 = \mathcal{M}_*^{(0)} \mathcal{L}_0 \mathcal{Z} \mathcal{J}_M \mathcal{J}_L = \mathcal{L}_*^{(0)} \mathcal{M}_0 \mathcal{Z} \mathcal{J}_L \mathcal{J}_M \\ &= \mathcal{L}_*^{(0)} \mathcal{L}_0 \mathcal{Z} \mathcal{J}_M^2 = \mathcal{L}_1 \mathcal{Z} \mathcal{J}_M^2. \end{aligned}$$

Inductively, we have

$$(3.11) \quad \mathcal{M}_k \mathcal{Z} \mathcal{J}_L^{2^k} = \mathcal{L}_k \mathcal{Z} \mathcal{J}_M^{2^k}$$

for any positive integer k . From the definitions of \mathcal{J}_M and \mathcal{J}_L in (3.5a) and (3.5b), respectively, it can be deduced that (3.11) can be rewritten as

$$(3.12) \quad \mathcal{M}_k \mathcal{Z} \begin{bmatrix} I & 0 \\ 0 & (J_s^H)^{2^k} \oplus I_\mu \end{bmatrix} = \mathcal{L}_k \mathcal{Z} \begin{bmatrix} J_s^{2^k} \oplus J_1^{2^k} & 0_\ell \oplus \Gamma_k \\ 0 & I_\ell \oplus J_1^{2^k} \end{bmatrix},$$

where

$$(3.13) \quad \Gamma_k = 2^{k-1} J_1^{2^{k-1}-1}$$

for any positive number k . From (3.6) and (3.13), we immediately obtain the following Lemma.

Lemma 3.1. *Let J_1 , J_s and Γ_k be defined in (3.6) and (3.13), respectively. Then Γ_k is invertible and satisfies*

$$(3.14) \quad \|\Gamma_k^{-1}\| = O(2^{-k}), \quad \|J_1^{2^k}\| = O(1), \quad \|J_s^{2^k} \oplus \Gamma_k^{-1}\| = O(2^{-k}).$$

Based on the Lemma 3.1, we shall prove the linear convergence of SDA under the condition **(A)**. We now partition \mathcal{Z} in (3.8) by

$$(3.15) \quad \mathcal{Z} = \begin{bmatrix} Z_1 & Z_3 \\ Z_2 & Z_4 \end{bmatrix},$$

where $Z_i \in \mathbb{R}^{n \times n}$, for $i = 1, 2, 3, 4$.

Theorem 3.2. *Suppose that the $(\mathcal{M}, \mathcal{L})$ in (1.5) satisfies the assumption **(A)** and that the DARE (1.1) has an almost stabilizing solution X_s with property **(P)**. Let $Z_{2b} = Z_2(:, 1 : \mu)$ and $Z_{4a} = Z_4(:, 1 : \ell)$. If the matrix $[Z_{4a} \ Z_{2b}] \in \mathbb{R}^{n \times n}$ is invertible, then Z_1 is invertible, $\widetilde{X}_s = Z_2 Z_1^{-1}$ is an almost stabilizing solution of DARE (2.8), and the sequences $\{A_k, G_k, H_k\}$ generated by Algorithm 2.1 satisfy*

- (1) $\limsup_{k \rightarrow \infty} \sqrt[k]{\|A_k\|} \leq \frac{1}{2}$.
- (2) $\limsup_{k \rightarrow \infty} \frac{\|H_{k+1} - \widetilde{X}_s\|}{\|H_k - \widetilde{X}_s\|} \leq \frac{1}{2}$, i.e., $H_k \rightarrow \widetilde{X}_s$ linearly with rate less than or equal to $\frac{1}{2}$. Moreover, $X_s = \widetilde{X}_s + Y$.

Proof. By Assumption **(A)** and assumptions in Lemma 1.1, $\text{span} \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix}$ and $\text{span} \begin{bmatrix} I \\ X \end{bmatrix}$ forms the same unique stable Lagrangian of $(\mathcal{M}_0, \mathcal{L}_0)$ corresponding to $J_s \oplus J_1$. Then Z_1 is invertible, and $\widetilde{X}_s = Z_2 Z_1^{-1}$ solves the DARE (2.8)[19].

On the other hand, substituting $(\mathcal{M}_k, \mathcal{L}_k)$ of (3.9) and \mathcal{Z} of (3.15) into (3.12), and comparing both sides we obtain

$$(3.16a) \quad A_k Z_1 = (Z_1 + G_k Z_2) (J_s^{2k} \oplus J_1^{2k}),$$

$$(3.16b) \quad A_k Z_3 \left((J_s^H)^{2k} \oplus I_\mu \right) = (Z_1 + G_k Z_2) (0_\ell \oplus \Gamma_k) \\ + (Z_3 + G_k Z_4) (I_\ell \oplus J_1^{2k}),$$

$$(3.16c) \quad -H_k Z_1 + Z_2 = A_k^\top Z_2 (J_s^{2k} \oplus J_1^{2k}),$$

$$(3.16d) \quad (-H_k Z_3 + Z_4) \left((J_s^H)^{2k} \oplus I_\mu \right) = A_k^\top Z_2 (0_\ell \oplus \Gamma_k) + A_k^\top Z_4 (I_\ell \oplus J_1^{2k}).$$

Postmultiplying (3.16c) by Z_1^{-1} and then substituting H_k into (3.16d), we have

$$(3.17) \quad (-\widetilde{X}_s Z_3 + A_k^\top Z_2 (J_s^{2k} \oplus J_1^{2k}) Z_1^{-1} Z_3 + Z_4) \left((J_s^H)^{2k} \oplus I_\mu \right) \\ = A_k^\top Z_2 (0_\ell \oplus \Gamma_k) + A_k^\top Z_4 (I_\ell \oplus J_1^{2k}).$$

Postmultiplying (3.17) by $(I_\ell \oplus \Gamma_k^{-1})$, we get

$$\begin{aligned}
 (3.18) \quad & A_k^\top \left[Z_2(0_\ell \oplus I_\mu) + Z_4(I_\ell \oplus J_1^{2^k} \Gamma_k^{-1}) \right. \\
 & \left. - Z_2(J_s^{2^k} \oplus J_1^{2^k}) Z_1^{-1} Z_3((J_s^H)^{2^k} \oplus \Gamma_k^{-1}) \right] \\
 & = (-\widetilde{X}_s Z_3 + Z_4)((J_s^H)^{2^k} \oplus \Gamma_k^{-1}).
 \end{aligned}$$

By Lemma 3.1 and assumptions in Theorem 3.2, for sufficient large k , the matrix $Z_2(0_\ell \oplus I_\mu) + Z_4(I_\ell \oplus J_1^{2^k} \Gamma_k^{-1})$ is invertible and the matrix $Z_2(J_s^{2^k} \oplus J_1^{2^k}) Z_1^{-1} Z_3((J_s^H)^{2^k} \oplus \Gamma_k^{-1})$ tends to 0. Therefore, $Z_2(0_\ell \oplus I_\mu) + Z_4(I_\ell \oplus J_1^{2^k} \Gamma_k^{-1}) - Z_2(J_s^{2^k} \oplus J_1^{2^k}) Z_1^{-1} Z_3((J_s^H)^{2^k} \oplus \Gamma_k^{-1})$ will be invertible for sufficiently large values of k . Then the sequence $\{A_k\}$ satisfies

$$(3.19) \quad \limsup_{k \rightarrow \infty} \sqrt[k]{\|A_k\|} \leq \limsup_{k \rightarrow \infty} \sqrt[k]{O(1)2^{-k}} = \frac{1}{2}.$$

From (3.16c), we get

$$\begin{aligned}
 & \limsup_{k \rightarrow \infty} \frac{\|H_{k+1} - \widetilde{X}_s\|}{\|H_k - \widetilde{X}_s\|} \leq \limsup_{k \rightarrow \infty} \sqrt[k]{\|H_k - \widetilde{X}_s\|} \\
 & = \limsup_{k \rightarrow \infty} \sqrt[k]{\|A_k^\top Z_2(J_s^{2^k} \oplus J_1^{2^k}) Z_1^{-1}\|} \leq \frac{1}{2}. \quad \blacksquare
 \end{aligned}$$

Corollary 3.3. *Assume that (A, B) is d -stabilizable and that the same conditions as in Theorem 3.2 hold. If the DARE (1.1) has a maximal solution X_+ , then it must coincide with the almost stabilizing solution X_s computed by SDA.*

Form Theorem 1.2, it can be seen that the maximal solution X_+ satisfies $R+B^\top X_+ B > 0$ and $\rho(A_F) \leq 1$. In addition, since (A, B) is d -stabilizable, the assumptions of Lemma 1.1 can be guaranteed. Therefore, the subspaces $\text{span} \begin{bmatrix} I \\ X_+ - Y \end{bmatrix}$ and $\text{span} \begin{bmatrix} I \\ \widetilde{X}_s \end{bmatrix}$ are unique stable Lagrangian subspaces of matrix pair $(\mathcal{M}_0, \mathcal{L}_0)$ corresponding to the spectrum of $J_s \oplus J_1$. Hence, this completes the proof. \blacksquare

Remark 3.2. *If the d -stabilizability of the pair (A, B) is replaced by a weaker condition (1.7), then the conclusion of Corollary 3.3 still holds.*

4. NUMERICAL EXAMPLES

The aim of this section is to illustrate the superior performance of the SDA, as compared to the Newton’s Method [8]. The flop counts for each iteration in

SDA and NTM is $\frac{23}{3}n^3$ and $30n^3$, respectively. We test some numerical examples satisfying Assumption(A) by the SDA and NTM. The convergence of NTM is guaranteed under the same conditions, and the rate of convergence is linear with ratio $\frac{1}{2}$.

Note that the NTM can be used to solve DARE with a more general case when R is singular, but we must assume that there is a symmetric solution X_+ of the inequality $\mathcal{R}(X_+) \geq 0$ for which $R + B^\top X_+ B > 0$. In the SDA, matrices G_0 and H_0 in (2.5) are only required to be symmetric. Under this relaxed condition, the existence of sequence $\{A_k, G_k, H_k\}$ is guaranteed. As mentioned before, the approximate solution computed by SDA algorithm is an almost stabilizing solution X_s of DARE (1.1) when all eigenvalues of A_F are in closed unit disk. The Newton's method converges to the maximal solution X_+ of DARE (1.1). In Corollary 3.3, we prove that these two solutions are coincident which is observed in our numerical experiments.

We report the numbers of iterations by "ITs", the CPU time by "CPU" for two algorithms, and the "Err" in SDA and NTM is defined by $\|H_{k+1} - H_k\|$ and $\|X_{k+1} - X_k\|$, respectively. We list five examples in this section. Example 4.1, 4.2 and 4.3 are identical to numerical examples in [5], which were presented originally in [13, 8, 21]. In the fourth example, the proven convergence rate of the SDA has been observed when the close loop matrix A_F has semi-simple eigenvalues on the unit circle. In the last example, we list the CPU time ratios of the SDA and NTM with increasing dimensions n .

For the residual of DARE, we use the "normalized" residual (DNRes) formula

$$\begin{aligned} & \text{DNRes} \\ & \equiv \frac{\| -\tilde{X} + A^\top \tilde{X} A + Q - (C + B^\top \tilde{X} A)^\top (R + B^\top \tilde{X} B)^{-1} (C + B^\top \tilde{X} A) \|}{\|\tilde{X}\| + \|A^\top \tilde{X} A\| + \|Q\| + \|(C + B^\top \tilde{X} A)^\top (R + B^\top \tilde{X} B)^{-1} (C + B^\top \tilde{X} A)\|}, \end{aligned}$$

proposed in [4], where \tilde{X} is an approximate solution to DARE.

All computations were performed in MATLAB/version 7.0 on a PC with a Intel Pentium-IV 3.2 GHZ processor and 2.5 GB main memory, using IEEE double-precision.

Example 4.1. [13]. For the following numerical data with singular A and R

$$\begin{aligned} A &= \begin{bmatrix} 0 & 1 \\ 0 & -1 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 0 \\ 2 & 1 \end{bmatrix}, \quad Q = \begin{bmatrix} -\frac{4}{11} & -\frac{4}{11} \\ -\frac{4}{11} & \frac{7}{11} \end{bmatrix} \\ C &= \begin{bmatrix} 3 & 1 \\ -1 & 7 \end{bmatrix}, \quad R = \begin{bmatrix} 9 & 3 \\ 3 & 1 \end{bmatrix}. \end{aligned}$$

Note that the rank of R is 1, i.e., R is singular. The symplectic pencil $(\mathcal{M}, \mathcal{L})$ has no eigenvalue on the unit circle. At each NTM iteration, Algorithm 2.2 solves a

Stein equation expensively [2]. From Table 1 we see that the SDA and NTM are both converge quadratically and the CPU time of the SDA is shorter than of NTM.

Table 1. Results for Example 4.1.

	SDA	NTM
DNRes	2.89×10^{-16}	1.6×10^{-16}
ITs	6	6
CPU	0.016	0.12
Err	1.1×10^{-10}	1.4×10^{-10}

Example 4.2. [8]. We consider the DARE with $n = m = 2$ defined by

$$A = \begin{bmatrix} 0 & -1 \\ 0 & 2 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}, \quad C = 0, \quad Q = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad R = \begin{bmatrix} 4 & 2 \\ 2 & 1 \end{bmatrix}.$$

We note that R is singular, and the symplectic pencil $(\mathcal{M}, \mathcal{L})$ has eigenvalues $\{0, 1, 1, \infty, \infty, \infty\}$. It can be easily seen $X = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$ which is the only solution of the DARE. The close-loop $A + BF$ has eigenvalues $\{0, 1\}$ and the elementary divisors of $\mathcal{M} - \lambda\mathcal{L}$ corresponding to the eigenvalue $\{1\}$ are of degree two. We can see that the convergence of the SDA and NTM are both linear with rate $\frac{1}{2}$. The numerical results are recorded in Table 2.

Table 2. Results for Example 4.2.

	SDA	NTM
DNRes	1.2×10^{-16}	1.6×10^{-16}
ITs	24	26
CPU	0.031	0.14
Err	3.0×10^{-8}	1.4×10^{-8}

Example 4.3. [21]. Consider the DARE (1.1) with

$$A = \begin{bmatrix} 0 & 10^{-1} & 0 \\ 0 & 0 & 10^{-1} \\ 0 & 0 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix},$$

and

$$Q = \begin{bmatrix} 10^5 & 0 & 0 \\ 0 & 10^3 & 0 \\ 0 & 0 & -10 \end{bmatrix}, \quad R = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}, \quad C = 0.$$

The exact solution of the DARE is $X = \text{diag}(10^5, 10^3, 0)$. Since the symplectic pair $(\mathcal{M}, \mathcal{L})$ has no eigenvalues on the unit circle, both methods converge quadratically.

Here, we choose the initial matrix $L_0 = 0$ so that $A - BL_0 = A$ is d-stable. From Table 3, it follows that the approximate solutions X from the SDA and NTM both have 16 significant digits.

Table 3. Results for Example 4.3.

	SDA	NTM
DNRes	4.6×10^{-16}	3.9×10^{-16}
ITs	2	2
CPU	0.031	0.047
Err	0	0

Example 4.4. Let r be an arbitrary number, and $R = \begin{bmatrix} 1 \\ r \end{bmatrix} \begin{bmatrix} 1 & r \end{bmatrix}$, $A = \begin{bmatrix} 2 + r^2 & 0 \\ 0 & 0 \end{bmatrix}$, $B = I_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$, $C = 0$, and $Q = I_2 - A^\top A + (C + B^\top A)^\top (R + B^\top B)^{-1} (C + B^\top A)$. It can be easily verified that the almost stabilizing solution $X = I_2$, and $A + BF = \begin{bmatrix} 1 & 0 \\ r & 0 \end{bmatrix}$ has eigenvalues 1, 0; i.e., assumption (A) holds. Newton's method needs to choose an initial matrix L_0 such that $A - BL_0$ is d-stable. It is easy to check that $L_0 = A$ satisfies the requirement. We list the absolute errors after the 20th iteration in the SDA and NTM in Table 4, the linear convergence rate with ratio $\frac{1}{2}$ can be observed.

Table 4. Results for Example 4.4.

ITs	Err(SDA)	Err(NTM)
20	1.26e-6	1.68e-6
21	6.29e-7	6.21e-7
22	3.15e-7	3.16e-7
23	1.58e-7	1.64e-7
24	8.01e-8	8.07e-8

Example 4.5. In this example, we run the algorithms on some randomly generated examples with the dimension n varying from 50 to 300. We shall construct $n \times n$ matrices A , B , Q and R such that the spectrum of $A + BF$ lies on the unit circle. In the first place, let U be a random unitary matrix and $A = 2U$. The solution X is a symmetric positive definite matrix, and R is a symmetric positive semidefinite matrix with one eigenvalue 0, $n - 1$ eigenvalues between 0 and 1. Let $B = X^{-\frac{1}{2}} \text{chol}(I - R)$, $C = \frac{1}{2} B^{-1} A - B^\top X A$, $Q = X - A^\top X A - A^\top B^{-2} A$. It is easy to check that the matrix $A + BF$ and the unitary matrix U are identical, as we have designed.

In Figure 4.1, we report a comparison of CPU times for the SDA and NTM for $n = 50, 100, 150, 200, 250, 300$. We also list the normalized residuals (NRes) in Table 5. From Table 5, we observe that the residuals of the SDA are smaller than those of NTM, up to 1-2 more digits of accuracy for all n . This indicates that the

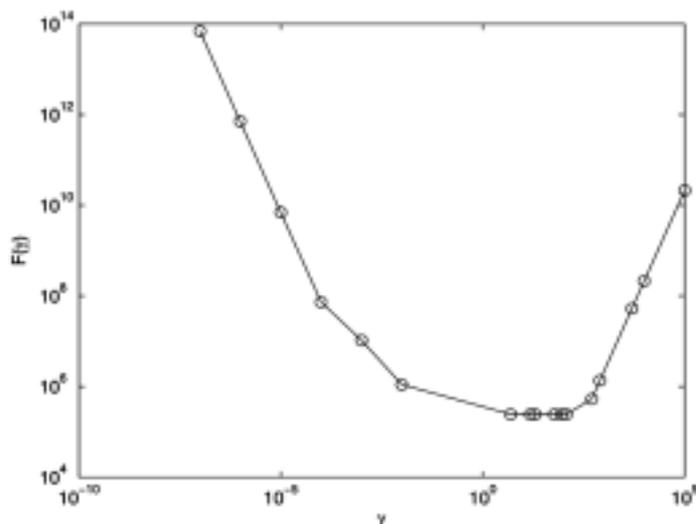


Fig. 4.1. The ratio of CPU times of Example 4.5.

Table 5. Results for Example 4.5.

Methods		NTM	SDA	Methods		NTM	SDA
$n = 50$	IT	17	19	$n=200$	IT	18	20
	CPU	0.41	0.17		CPU	14	2.1
	NRes	4.6e-12	2.3e-13		NRes	5.2e-14	9.6-14
$n = 100$	IT	18	19	$n=250$	IT	17	21
	CPU	2.3	0.33		CPU	29	4.9
	NRes	4.3e-12	6.1e-13		NRes	4.1e-12	5.6e-14
$n = 150$	IT	18	20	$n=300$	IT	17	20
	CPU	6.8	1.1		CPU	57	6.4
	NRes	5.4e-12	1.6e-13		NRes	2.3e-12	7.8e-14

SDA computes more accurate solutions than NTM, generally. In Figure 4.1, we see that the CPU time of the SDA is approximately 10% to 30% of that of NTM.

5. CONCLUDING REMARKS

In this paper, we propose the structured doubling algorithm for finding the stabilizing or almost stabilizing solution of DARE (1.1). In Theorem 3.1 and 3.2,

we prove quadratic and global linear with ratio $\frac{1}{2}$ convergence for SDA algorithms, respectively. The convergence behavior is similar to that of Newton's method. We prove in Corollary 3.3 that the almost stabilizing solution computed by SDA is the same as the maximal solution computed by Newton's method. However, in each Newton's iteration, a Stein equation must be solved, which is rather expensive. Numerical examples show that our structured doubling algorithm is efficient, outperforming Newton's method.

ACKNOWLEDGMENT

We would like to thank Professor Eric King-Wah Chu from Monash University for many helpful suggestions and valuable comments.

REFERENCES

1. Z. Bai, J. Demmel and M. Gu, An inverse free parallel spectral divide and conquer algorithm for nonsymmetric eigenproblem, *Numer. Math.*, **76** (1997), 279-308.
2. R. H. Bartels and G. W. Stewart, Solution of the matrix equation $AX + XB = C$, *Comm. ACM.*, **15** (1972), 820-826.
3. M. S. Bazaraa, H. D. Sherali and C. M. Shetty, *Nonlinear Programming*, John Wiley, 1993.
4. P. Benner and R. Byers, Evaluating Products of matrix pencils and collapsing matrix products, *Num. Lin. Alg. Appl.*, **8** (2001), 357-380.
5. P. Benner, A. J. Laub and V. Mehrmann, *A collection of benchmark examples for the numerical solution of algebraic Riccati equations II: discrete-time case*, Technical Report SPC 95-23, Fakultät für Mathematik, Technische Universität Chemnitz-Zwickau, FRG, 1995.
6. E. K.-W. Chu, H.-Y. Fan and W.-W. Lin, Structure-preserving algorithms for periodic discrete-time algebraic Riccati equations, *Int. J. Control*, **77** (2004), 767-788.
7. A. Emami-Naeini and G. Franklin, Comments on the numerical solution of the discrete-time algebraic Riccati equation, *IEEE Trans. Automat. Control*, **25** (1980), 1015-1016.
8. C.-H. Guo, Newton's method for discrete algebraic Riccati equations with the closed-loop matrix has eigenvalues on the unit circle, *SIAM J. Matrix Anal. Appl.*, **20** (1998), 279-294.
9. N. J. Higham, Stability of the diagonal pivoting method with partial pivoting, *SIAM J. Matrix Anal. Appl.*, **18** (1997), 52-65.
10. N. J. Higham, *Accuracy and Stability of Numerical Algorithms*, 2nd ed., SIAM, 2002.
11. T.-W. Hwang, E. K.-W. Chu and W.-W. Lin, A generalized structure-preserving doubling algorithm for generalized discrete-time algebraic Riccati equations, *Int. J. Control*, **78** (2005), 1063-1075.

12. T.-W. Hwang and W.-W. Lin, Structured doubling algorithms for weakly stabilizing Hermitian solutions of algebraic Riccati equations, *Lin. Alg. Appl.*, **430** (2009), 1452-1478.
13. V. Ionescu and M. Weiss, On computing the stabilizing solution of the discrete-time Riccati equation, *Lin. Alg. Appl.* **174** (1992), 229-238.
14. P. Lancaster and L. Rodman, *Algebraic Riccati Equations*, Clarendon Press, Oxford, 1995.
15. W.-W. Lin, V. Mehrmann and H. Xu, Canonical forms for Hamiltonian and Symplectic matrices and pencils, *Lin. Alg. Appl.*, **302/303** (1999), 469-533.
16. W.-W. Lin and S.-F. Xu, Convergence analysis of structure-preserving doubling algorithms for Riccati-type matrix equations, *SIAM Matrix Anal. Appl.*, **28** (2006), 26-39.
17. A. N. Malyshev, Parallel algorithm for solving some spectral problems of linear algebra, *Lin. Alg. Appl.*, **188/189** (1993), 489-520.
18. MathWorks, *MATLAB User's Guide (for Windows version)*, The Math Works, Inc., 2002.
19. V. Mehrmann, *The Autonomous Linear Quadratic Control Problem*, Springer-Verlag, 1991.
20. C. Paige and C. Van Loan, A Schur decomposition for Hamiltonian matrices, *Lin. Alg. Appl.*, **41** (1981), 11-32.
21. J. G. Sun, Sensitivity of the discrete-time algebraic Riccati equation, *Lin. Alg. Appl.*, **255-276** (1998), 595-615.
22. P. Van Dooren, A generalized eigenvalue approach for solving Riccati equations, *SIAM J. Sci. Comput.*, **2** (1981), 121-135.

Chun-Yueh Chiang
Center for General Education,
National Formosa University,
Huwei 632, Taiwan
E-mail: chiang@nfu.edu.tw

Hung-Yuan Fan
Department of Mathematics,
National Taiwan Normal University,
Taipei 116, Taiwan
E-mail: hyfan@math.ntnu.edu.tw

Wen-Wei Lin
Department of Applied Mathematics,
National Chiao Tung University,
Hsinchu 300, Taiwan
E-mail: wwlin@math.nctu.edu.tw