

## EXPLICIT SOLUTIONS FOR A RICCATI EQUATION FROM TRANSPORT THEORY\*

VOLKER MEHRMANN<sup>†</sup> AND HONGGUO XU<sup>‡</sup>

*In memoriam of Gene H. Golub*

**Abstract.** We derive formulas for the minimal positive solution of a particular nonsymmetric Riccati equation arising in transport theory. The formulas are based on the eigenvalues of an associated matrix. We use the formulas to explore some new properties of the minimal positive solution and to derive fast and highly accurate numerical methods. Some numerical tests demonstrate the properties of the new methods.

**Key words.** nonsymmetric Riccati equation, secular equation, eigenvalues, minimal positive solution, Cauchy matrix, transport theory, quadrature formula

**AMS subject classifications.** 15A24, 65F15, 82C70, 65H05

**DOI.** 10.1137/070708743

**1. Introduction.** We consider nonsymmetric matrix Riccati equations of the special form

$$(1.1) \quad XA + DX - XBX - C = 0,$$

with

$$A = \Gamma - pe^T, \quad D = \Delta - ep^T, \quad B = pp^T, \quad C = ee^T,$$

where

$$\Gamma := \text{diag}(\gamma_1, \dots, \gamma_n), \quad \Delta := \text{diag}(\delta_1, \dots, \delta_n),$$

$$p = [p_1, \dots, p_n]^T, \quad e = [1, \dots, 1]^T,$$

and  $\gamma_n > \dots > \gamma_1 > 0$ ,  $\delta_n > \dots > \delta_1 > 0$ , and  $p_1, \dots, p_n > 0$ .

Such Riccati equations arise in Markov models [28] and in nuclear physics [8, 18, 22]. In the latter application, to study the transport of particles, one introduces integral equations of the form

$$(1.2) \quad \left[ \frac{1}{x + \alpha} + \frac{1}{y - \alpha} \right] T(x, y) = \beta \left[ 1 + \frac{1}{2} \int_{-\alpha}^1 \frac{T(t, y)}{t + \alpha} dt \right] \left[ 1 + \frac{1}{2} \int_{\alpha}^1 \frac{T(x, t)}{t - \alpha} dt \right],$$

where the unknown function  $T(x, y) : [-\alpha, 1] \times [\alpha, 1] \mapsto \mathbb{R}^+$  is called the *scattering function*,  $\alpha \in [0, 1)$  is an angular shift, and  $\beta \in [0, 1]$  is the average of the total

---

\*Received by the editors November 20, 2007; accepted for publication (in revised form) by J. H. Brandts June 4, 2008; published electronically October 16, 2008.

<http://www.siam.org/journals/simax/30-4/70874.html>

<sup>†</sup>Institut für Mathematik, TU Berlin, Str. des 17. Juni 136, D-10623 Berlin, Germany (mehrman@math.tu-berlin.de). This author's research was partially supported by Deutsche Forschungsgemeinschaft, through the DFG Research Center MATHEON Mathematics for Key Technologies in Berlin.

<sup>‡</sup>Department of Mathematics, University of Kansas, Lawrence, KS 44045 (xu@math.ku.edu). This author's research was partially supported by the University of Kansas General Research Fund allocation # 2301717 and by Deutsche Forschungsgemeinschaft through the DFG Research Center MATHEON Mathematics for Key Technologies in Berlin.

number of particles emerging from a collision. (Here  $\mathbb{R}^+$  denotes the set of positive real numbers.)

To solve this integral equation numerically, one approximates the integrals via classical quadrature formulas [29]. For this the function  $T(x, y)$  is approximated via a matrix  $X = [x_{ij}]$ , where  $x_{ij}$  is an approximation of  $T(\mu_i, \nu_j)$  with  $\mu_i, \nu_j$  being the  $i$ th and  $j$ th nodes of the quadrature formula on  $[-\alpha, 1]$  and  $[\alpha, 1]$ , respectively; see, e.g., [18].

In this discretization the matrix  $X$  has to satisfy the matrix Riccati equation (1.1) with coefficient matrices

$$(1.3) \quad \gamma_j = \frac{1}{\beta(1-\alpha)\omega_j}, \quad \delta_j = \frac{1}{\beta(1+\alpha)\omega_j}, \quad p_j = \frac{c_j}{2\omega_j},$$

for  $j = 1, 2, \dots, n$ , where  $\{c_j\}_{j=1}^n, \{\omega_j\}_{j=1}^n$  are the sets of weights and nodes of the specific quadrature rule that is used on the interval  $[0, 1]$ . These typically satisfy

$$(1.4) \quad c_1, \dots, c_n > 0, \quad \sum_{j=1}^n c_j = 1; \quad 1 > \omega_1 > \dots > \omega_n > 0.$$

In [20] it is shown that the Riccati equation (1.1) has two entrywise positive solutions  $X = [x_{ij}], Y = [y_{ij}] \in \mathbb{R}^{n,n}$ , which satisfy  $X \leq Y$ , where we use the notation that  $X \leq Y$  if  $x_{ij} \leq y_{ij}$  for all  $i, j = 1, \dots, n$ .

In the applications from transport theory, only  $X$ , the smaller of the two positive solutions, is of interest. Therefore, in this paper we consider only the computation of the minimal positive solution  $X$ . The computation of this minimal solution has been investigated in several publications. Various direct and iterative methods [1, 2, 12, 13, 14, 15, 16, 17, 18, 19, 25] have been proposed by either directly solving the Riccati equation or by computing specific invariant subspaces of the  $2n \times 2n$  matrix

$$(1.5) \quad H = \begin{bmatrix} A & -B \\ C & -D \end{bmatrix}$$

that is formed from the coefficient matrices.

In [20] even an explicit solution formula has been derived that is based on the eigenvalues  $H$ . Motivated by this result, we derive different explicit formulas, one of which is mathematically equivalent to the one in [20], but of a much simpler form. We will use these formulas to derive both entrywise and normwise bounds for the solution matrix and show that the entries of the solution have a graded entry property. We will also use the formulas to develop fast and highly accurate numerical algorithms for the minimal positive solution of (1.1).

This paper is organized as follows. In section 2, we will reformulate the associated eigenvalue problem via an appropriate balancing strategy. We use the associated secular function to derive some properties of the eigenvalues of  $H$ . In section 3, we then derive four formulas for the minimal positive solution based on the eigenvalues. Entrywise and normwise bounds for the minimal positive solution are provided in section 4. Numerical algorithms and an error analysis are presented in section 5 and some numerical examples are shown in section 6. A conclusion is given in section 7.

Throughout this paper,  $\lambda(A)$  denotes the spectrum of a square matrix  $A$ , and  $I_n$  (or simply  $I$ ) is the  $n \times n$  identity matrix. The norm used in this paper is the spectral norm.

**2. Spectral properties of the matrix  $H$ .** In this section we analyze the spectral properties of the matrix  $H$  in (1.5) defined by the coefficient matrices of (1.1).

In order for all of the eigenvalues of  $H$  to be real, we assume that the condition

$$(2.1) \quad 1 - \sum_{j=1}^n p_j \left( \frac{1}{\gamma_j} + \frac{1}{\delta_j} \right) \geq 0$$

holds. The transport problem with the coefficients defined in (1.3) and (1.4) is a special case where this assumption is satisfied.

The first step in our analysis is a balancing of the coefficient matrices. Since the entries of the vector  $p$  are positive, we may define

$$\Phi := \text{diag}(\sqrt{p_1}, \dots, \sqrt{p_n}), \quad \phi := [\sqrt{p_1}, \dots, \sqrt{p_n}]^T.$$

Using  $\Phi$  to scale the Riccati equation (1.1) via

$$\begin{aligned} \tilde{X} &= \Phi X \Phi, \\ \tilde{A} &= \Phi^{-1} A \Phi = \Gamma - \phi \phi^T, \\ \tilde{D} &= \Phi D \Phi^{-1} = \Delta - \phi \phi^T, \\ \tilde{B} &= \Phi^{-1} B \Phi^{-1} = \phi \phi^T, \\ \tilde{C} &= \Phi C \Phi = \phi \phi^T = \tilde{B}, \end{aligned}$$

we obtain the equivalent Riccati equation

$$(2.2) \quad \tilde{X} \tilde{A} + \tilde{D} \tilde{X} - \tilde{X} \tilde{B} \tilde{X} - \tilde{B} = 0,$$

and obviously,  $X$  is a solution to (1.1) if and only if  $\tilde{X} = \Phi X \Phi$  is a solution to (2.2).

For the associated matrix formed from the coefficients we then have

$$(2.3) \quad \begin{aligned} \tilde{H} &= \begin{bmatrix} \Phi^{-1} & 0 \\ 0 & \Phi \end{bmatrix} H \begin{bmatrix} \Phi & 0 \\ 0 & \Phi^{-1} \end{bmatrix} \\ &= \begin{bmatrix} \tilde{A} & -\tilde{B} \\ \tilde{B} & -\tilde{D} \end{bmatrix} = \begin{bmatrix} \Gamma & 0 \\ 0 & -\Delta \end{bmatrix} - \begin{bmatrix} \phi \\ -\phi \end{bmatrix} \begin{bmatrix} \phi \\ \phi \end{bmatrix}^T, \end{aligned}$$

and we see that  $\tilde{H}$  is similar to  $H$  and is a rank-one modification of a diagonal matrix, which is analogous to the real symmetric rank-one updating problem discussed in [9]. It follows that the eigenvalues of  $\tilde{H}$  can be obtained cheaply and accurately via the solution of secular equations by using a method similar to the one discussed in [10, section 8.5].

Furthermore, it is well known (see, e.g., [23]) that  $\tilde{X}$  is a solution to (2.2) if and only if  $\tilde{X}$  satisfies the invariant subspace equation

$$\tilde{H} \begin{bmatrix} I \\ \tilde{X} \end{bmatrix} = \begin{bmatrix} I \\ \tilde{X} \end{bmatrix} (\tilde{A} - \tilde{B} \tilde{X}).$$

In [20] it was shown (for the original solution  $X$ ) that  $\tilde{X}$  is the minimal positive solution if and only if all of the eigenvalues of  $\tilde{A} - \tilde{B} \tilde{X}$  are nonnegative.

In order to analyze the properties of the matrix  $\tilde{H}$  and thus also of the similar matrix  $H$ , we first derive some properties of the eigenvalues of  $\tilde{H}$ .

Consider the rational function

$$(2.4) \quad \chi(\lambda) = 1 + \sum_{j=1}^n \frac{p_j}{\lambda - \gamma_j} - \sum_{j=1}^n \frac{p_j}{\lambda + \delta_j}.$$

Then, since

$$(2.5) \quad \det(\lambda I - \tilde{H}) = \chi(\lambda) \left( \prod_{j=1}^n (\lambda - \gamma_j)(\lambda + \delta_j) \right),$$

it follows that the eigenvalues of  $\tilde{H}$  are just the roots of the *secular equation*  $\chi(\lambda) = 0$ , and thus the computation of the spectrum of  $\tilde{H}$  can be carried out very efficiently by solving the secular equation; see [11, 27]. Furthermore, we have the following interlacing properties.

**LEMMA 2.1.** *Consider the matrix  $\tilde{H}$  defined via the coefficients of the Riccati equation (2.2), and suppose that (2.1) holds. Then  $\tilde{H}$  has  $2n$  real eigenvalues,  $-\nu_n < \dots < -\nu_1 \leq 0$ ,  $0 \leq \lambda_1 < \dots < \lambda_n$ , that satisfy the inequalities*

$$0 \leq \nu_1 < \delta_1 < \nu_2 < \delta_2 < \dots < \nu_{n-1} < \delta_{n-1} < \nu_n < \delta_n$$

and

$$0 \leq \lambda_1 < \gamma_1 < \lambda_2 < \gamma_2 < \dots < \lambda_{n-1} < \gamma_{n-1} < \lambda_n < \gamma_n.$$

Moreover, the following cases can be considered:

1.  $\nu_1 = 0$  and  $\lambda_1 > 0$  if and only if  $\chi(0) = 0$  and  $\chi'(0) > 0$ .
2.  $\nu_1 > 0$  and  $\lambda_1 = 0$  if and only if  $\chi(0) = 0$  and  $\chi'(0) < 0$ .
3.  $\nu_1 = \lambda_1 = 0$  if and only if  $\chi(0) = \chi'(0) = 0$ . In this case,  $\tilde{H}$  has a  $2 \times 2$  Jordan block associated with the eigenvalue 0.

*Proof.* The proof is basically given already in [20] based on the properties of the secular function  $\chi(\lambda)$ . Note that assumption (2.1) implies that  $\chi(0) \geq 0$ .

The second part of the third case has already been shown in [12] in a more general setting.  $\square$

**Remark 2.2.** Suppose the quadrature formula that is used to discretize the integral equation (1.2) is of order greater than or equal to 3; i.e.,

$$\sum_{j=1}^n c_j w_j^k = \frac{1}{k+1}, \quad k = 0, 1, 2, 3.$$

With (1.3) it is easily verified that

$$\begin{aligned} \chi(0) &= 1 - \sum_{j=1}^n \left( \frac{p_j}{\gamma_j} + \frac{p_j}{\delta_j} \right) = 1 - \beta \sum_{j=1}^n c_j = 1 - \beta, \\ \chi'(0) &= \sum_{j=1}^n \left( -\frac{p_j}{\gamma_j^2} + \frac{p_j}{\delta_j^2} \right) = 2\alpha\beta^2 \sum_{j=1}^n c_j w_j = \alpha\beta^2, \\ \chi''(0) &= -2 \sum_{j=1}^n \left( \frac{p_j}{\gamma_j^3} + \frac{p_j}{\delta_j^3} \right) = -2(1 + 3\alpha^2)\beta^3 \sum_{j=1}^n c_j w_j^2 = -\frac{2}{3}(1 + 3\alpha^2)\beta^3, \\ \chi'''(0) &= 6 \sum_{j=1}^n \left( -\frac{p_j}{\gamma_j^4} + \frac{p_j}{\delta_j^4} \right) = 24\alpha(1 + \alpha^2)\beta^4 \sum_{j=1}^n c_j w_j^3 = 6\alpha(1 + \alpha^2)\beta^4. \end{aligned}$$

Since  $\chi'(0) \geq 0$ , we have that case 1 in Lemma 2.1 happens when  $\beta = 1$  and  $\alpha > 0$  and case 3 happens when  $\beta = 1$  and  $\alpha = 0$ . Case 2 will never happen.

**3. Formulas for the minimal positive solution.** In this section we will derive explicit formulas for the minimal positive solution of (1.1) in terms of the eigenvalues  $-\nu_1, \dots, -\nu_n, \lambda_1, \dots, \lambda_n$  of  $H$  (or  $\tilde{H}$ ). For this we need the following lemma.

LEMMA 3.1. *Suppose in the following that  $\tilde{X} \in \mathbb{R}^{n,n}$ . The following statements are equivalent.*

- (a)  $\tilde{X}$  is the minimal positive solution of (2.2).
- (b)  $\tilde{X}$  satisfies

$$\tilde{H} \begin{bmatrix} I_n \\ \tilde{X} \end{bmatrix} = \begin{bmatrix} I_n \\ \tilde{X} \end{bmatrix} \tilde{R}_1,$$

where  $\tilde{R}_1 = \tilde{A} - \tilde{B}\tilde{X}$  and  $\sigma(\tilde{R}_1) = \{\lambda_1, \dots, \lambda_n\}$ .

- (c)  $\tilde{X}^T$  is the minimal positive solution to the dual Riccati equation

$$(3.1) \quad \tilde{Y}\tilde{D} + \tilde{A}\tilde{Y} - \tilde{Y}\tilde{B}\tilde{Y} - \tilde{B} = 0.$$

- (d)  $\tilde{X}$  satisfies

$$(3.2) \quad \tilde{H} \begin{bmatrix} \tilde{X}^T \\ I_n \end{bmatrix} = \begin{bmatrix} \tilde{X}^T \\ I_n \end{bmatrix} \tilde{R}_2,$$

where  $\tilde{R}_2 = -(\tilde{D} - \tilde{B}\tilde{X}^T)$  and  $\sigma(\tilde{R}_2) = \{-\nu_1, \dots, -\nu_n\}$ .

*Proof.* The equivalence of (a) and (b) is given in [20]. The equivalence between (a) and (c) is obvious by taking the transpose on both sides of (2.2) or (3.1). The equivalence between (c) and (d) is shown in [12].  $\square$

With formulas for  $\tilde{R}_1, \tilde{R}_2$  as in Lemma 3.1 and the formulas for  $\tilde{A}, \tilde{D}$  and  $\tilde{B}$ , it follows that the minimal positive solution  $\tilde{X}$  of (2.2) satisfies the following relations:

$$(3.3) \quad \Gamma - \phi\tilde{\xi}^T = \tilde{R}_1, \quad \sigma(\tilde{R}_1) = \{\lambda_1, \dots, \lambda_n\},$$

$$(3.4) \quad \Delta - \phi\tilde{\eta}^T = -\tilde{R}_2, \quad \sigma(-\tilde{R}_2) = \{\nu_1, \dots, \nu_n\},$$

$$(3.5) \quad \tilde{X}\Gamma + \Delta\tilde{X} = \tilde{\eta}\tilde{\xi}^T,$$

where

$$\tilde{\xi} = (I + \tilde{X}^T)\phi, \quad \tilde{\eta} = (I + \tilde{X})\phi.$$

The last equation is a reformulation of (2.2). It thus follows that if the vectors  $\tilde{\xi}$  and  $\tilde{\eta}$  can be determined, then  $\tilde{X}$  can be easily formulated based on the simple Sylvester equation (3.5).

The following result shows that  $\tilde{\xi}$  and  $\tilde{\eta}$  can be determined based on the relations (3.3) and (3.4).

**PROPOSITION 3.2** (see [26]). *Suppose that matrices  $A, B$  are given such that  $A = \text{diag}(a_1, \dots, a_n)$  with distinct diagonal entries  $a_1, \dots, a_n \in \mathbb{R}$ , and  $B \in \mathbb{R}^{n,n}$  with  $\lambda(B) = \{b_1, \dots, b_n\}$  for distinct  $b_1, \dots, b_n \in \mathbb{R}$ .*

*Let  $q_1, q_2, \dots, q_n \in \mathbb{R} \setminus \{0\}$  and define*

$$q = [q_1, q_2, \dots, q_n]^T, \quad Q = \text{diag}(q_1, q_2, \dots, q_n)$$

as well as

$$f = \left[ \frac{\prod_{j=1}^n (a_1 - b_j)}{\prod_{j \neq 1} (a_1 - a_j)}, \dots, \frac{\prod_{j=1}^n (a_k - b_j)}{\prod_{j \neq k} (a_k - a_j)}, \dots, \frac{\prod_{j=1}^n (a_n - b_j)}{\prod_{j \neq n} (a_n - a_j)} \right]^T.$$

If a vector  $z \in \mathbb{R}^n$  satisfies  $A - qz^T = B$ , then

$$(3.6) \quad z = Q^{-1}f = \left[ \frac{f_1}{q_1}, \dots, \frac{f_n}{q_n} \right]^T.$$

Using (3.6), (3.3), (3.4), and (3.5), we obtain the following explicit formulas for  $X$ .

**THEOREM 3.3.** *Consider the Riccati equation (1.1). Introduce for  $k = 1, \dots, n$  the scalar quantities*

$$\xi_k = \frac{\prod_{j=1}^n (\gamma_k - \lambda_j)}{\prod_{j \neq k} (\gamma_k - \gamma_j)}, \quad \eta_k = \frac{\prod_{j=1}^n (\delta_k - \nu_j)}{\prod_{j \neq k} (\delta_k - \delta_j)}, \quad \kappa_k = \frac{\prod_{j=1}^n (\gamma_k + \delta_j)}{\prod_{j=1}^n (\gamma_k + \nu_j)}, \quad \epsilon_k = \frac{\prod_{j=1}^n (\delta_k + \gamma_j)}{\prod_{j=1}^n (\delta_k + \lambda_j)},$$

the associated vectors and matrices

$$(3.7) \quad \begin{aligned} \xi &= [\xi_1, \dots, \xi_n]^T, & \Xi &= \text{diag}(\xi_1, \dots, \xi_n), \\ \eta &= [\eta_1, \dots, \eta_n]^T, & E &= \text{diag}(\eta_1, \dots, \eta_n), \\ \kappa &= [\kappa_1, \dots, \kappa_n]^T, & K &= \text{diag}(\kappa_1, \dots, \kappa_n), \\ \epsilon &= [\epsilon_1, \dots, \epsilon_n]^T, & \mathcal{E} &= \text{diag}(\epsilon_1, \dots, \epsilon_n), \end{aligned}$$

and the Cauchy matrix

$$\Theta = \left[ \frac{1}{\delta_i + \gamma_j} \right].$$

Let

$$P = \text{diag}(p_1, \dots, p_n),$$

with the  $p_i$  defined in (1.1). Then we have the following solution formulas for (1.1):

$$(3.8) \quad X = P^{-1}E\Theta\xi P^{-1},$$

$$(3.9) \quad X = P^{-1}E\Theta K,$$

$$(3.10) \quad X = \mathcal{E}\Theta\xi P^{-1},$$

$$(3.11) \quad X = \mathcal{E}\Theta K.$$

*Proof.* To prove the formulas, we apply Proposition 3.2 to (3.3) and obtain

$$\tilde{\xi} = \Phi^{-1}\xi,$$

where  $\xi$  is defined in (3.7). Similarly, from (3.4) we obtain

$$\tilde{\eta} = \Phi^{-1}\eta,$$

where  $\eta$  is defined in (3.7). By solving the Sylvester equation (3.4) we obtain

$$\tilde{X} = \Phi^{-1}E\Theta\xi\Phi^{-1},$$

with  $E, \xi$  as in (3.7). Then, (3.8) follows by using  $X = \Phi^{-1}\tilde{X}\Phi^{-1}$  and  $P = \Phi^2$ .

In order to get the other formulas we need only show that  $\xi = PK$  and  $E = P\mathcal{E}$ .

Since  $-\nu_1, \dots, -\nu_n, \lambda_1, \dots, \lambda_n$  are the eigenvalues of  $\tilde{H}$ , it follows from (2.5) that

$$(3.12) \quad \prod_{j=1}^n (\lambda - \lambda_j) \prod_{j=1}^n (\lambda + \nu_j) = \sum_{m=1}^n p_m \prod_{j \neq m} (\lambda - \gamma_j) \prod_{j=1}^n (\lambda + \delta_j) - \sum_{m=1}^n p_m \prod_{j=1}^n (\lambda - \gamma_j) \prod_{j \neq m} (\lambda + \delta_j) + \prod_{j=1}^n (\lambda - \gamma_j) \prod_{j=1}^n (\lambda + \delta_j).$$

By inserting  $\lambda = \gamma_k$ , we obtain

$$\prod_{j=1}^n (\gamma_k - \lambda_j) \prod_{j=1}^n (\gamma_k + \nu_j) = p_k \prod_{j \neq k} (\gamma_k - \gamma_j) \prod_{j=1}^n (\gamma_k + \delta_j),$$

which implies that

$$\xi_k = p_k \kappa_k, \quad k = 1, 2, \dots, n.$$

We then have  $\xi = PK$ .

Similarly, by inserting  $\lambda = -\delta_k$  in (3.12) we get

$$\eta_k = p_k \epsilon_k, \quad k = 1, \dots, n,$$

and thus  $E = P\mathcal{E}$ . Then the other formulas follow.  $\square$

Note that formula (3.9) needs only the eigenvalues  $-\nu_1, \dots, -\nu_n$ , while formula (3.10) needs only the eigenvalues  $\lambda_1, \dots, \lambda_n$ . Numerically, these two formulas provide very cheap procedures to compute the minimal solution  $X$  of (1.1).

*Remark 3.4.* In [20] an explicit formula for the minimal solution of (1.1) was already given that is equivalent to (3.10). However, there a different expression for  $\epsilon_k$  was introduced as

$$\epsilon_k = 1 + \sum_{m=1}^n \frac{1}{\delta_k + \lambda_m} \frac{\prod_{j=1}^n (\gamma_j - \lambda_m)}{\prod_{j \neq m} (\lambda_j - \lambda_m)}.$$

This expression is less compact and its evaluation has a higher complexity than the expression in Theorem 3.3.

In this section we have derived new explicit formulas for the minimal solution  $X$  of (1.1) and we will use them in the next section to derive some further properties of  $X$ .

**4. Properties and bounds for the minimal positive solution.** The simple expressions of the quantities  $\xi_k, \kappa_k, \eta_k, \epsilon_k$  in the explicit formulas (3.8)–(3.11) and the eigenvalue interlacing property for the eigenvalues of  $\tilde{H}$  allow one to derive further properties of the minimal positive solution of (1.1). For this we first prove the following lemma.

LEMMA 4.1. *The coefficients  $\gamma_k, \delta_k$  in (1.1), the eigenvalues  $-\nu_k, \lambda_k$  of  $\tilde{H}$  in (2.3), and the quantities  $\xi_k, \eta_k, \kappa_k, \epsilon_k, k = 1, \dots, n$  in (3.7) satisfy the following inequalities.*

1.

$$0 < a_k < \eta_k < \delta_k - \nu_1 \leq \delta_k, \quad 0 < b_k < \xi_k < \gamma_k - \lambda_1 \leq \gamma_k,$$

$$1 < \epsilon_k < \frac{\delta_k + \gamma_n}{\delta_k + \lambda_1} \leq \frac{\delta_k + \gamma_n}{\delta_k}, \quad 1 < \kappa_k < \frac{\gamma_k + \delta_n}{\gamma_k + \nu_1} \leq \frac{\gamma_k + \delta_n}{\gamma_k},$$

where

$$a_k = \begin{cases} \frac{(\delta_k - \nu_k)(\nu_{k+1} - \delta_k)}{\delta_n - \delta_k}, & 1 \leq k < n, \\ \delta_n - \nu_n, & k = n, \end{cases}$$

$$b_k = \begin{cases} \frac{(\gamma_k - \lambda_k)(\lambda_{k+1} - \gamma_k)}{\gamma_n - \gamma_k}, & 1 \leq k < n, \\ \gamma_n - \lambda_n, & k = n. \end{cases}$$

2.

$$1 < \epsilon_n < \epsilon_{n-1} < \dots < \epsilon_1, \quad 1 < \kappa_n < \kappa_{n-1} < \dots < \kappa_1.$$

*Proof.* To prove the first part, we use the interlacing property in Lemma 2.1 and obtain

$$0 < \frac{\delta_k - \nu_j}{\delta_k - \delta_{j-1}} < 1, \quad 1 < j \leq k; \quad \frac{\delta_k - \nu_j}{\delta_k - \delta_j} > 1, \quad 1 \leq j < k,$$

and

$$0 < \frac{\delta_k - \nu_j}{\delta_k - \delta_j} < 1, \quad k < j \leq n; \quad \frac{\delta_k - \nu_{j+1}}{\delta_k - \delta_j} > 1, \quad k < j < n.$$

For  $1 \leq k < n$

$$\eta_k = \frac{(\delta_k - \nu_k)(\delta_k - \nu_{k+1})}{\delta_k - \delta_n} \prod_{j=1}^{k-1} \frac{\delta_k - \nu_j}{\delta_k - \delta_j} \prod_{j=k+1}^{n-1} \frac{\delta_k - \nu_{j+1}}{\delta_k - \delta_j} > a_k,$$



and

$$\eta_k = (\delta_k - \nu_1) \prod_{j=1}^{k-1} \frac{\delta_k - \nu_{j+1}}{\delta_k - \delta_j} \prod_{j=k+1}^n \frac{\delta_k - \nu_j}{\delta_k - \delta_j} < \delta_k - \nu_1 \leq \delta_k.$$

Finally, for  $k = n$  we obtain

$$\eta_n = (\delta_n - \nu_n) \prod_{j=1}^{n-1} \frac{\delta_n - \nu_j}{\delta_n - \delta_j} > \delta_n - \nu_n =: a_n$$

and

$$\eta_n = (\delta_n - \nu_1) \prod_{j=1}^{n-1} \frac{\delta_n - \nu_{j+1}}{\delta_n - \delta_j} < \delta_n - \nu_1 \leq \delta_n.$$

This proves the inequalities for the  $\eta_k$ , and clearly we have  $a_k > 0$  for  $k = 1, \dots, n$ .

The inequalities for the  $\xi_k$  can be derived in the same way by using the interlacing property for the eigenvalues  $\lambda_1, \dots, \lambda_n$ . This interlacing property also gives

$$\epsilon_k = \prod_{j=1}^n \frac{\delta_k + \gamma_j}{\delta_k + \lambda_j} > 1$$

and

$$\epsilon_k = \frac{\delta_k + \gamma_n}{\delta_k + \lambda_1} \prod_{j=1}^{n-1} \frac{\delta_k + \gamma_j}{\delta_k + \lambda_{j+1}} < \frac{\delta_k + \gamma_n}{\delta_k + \lambda_1} \leq \frac{\delta_k + \gamma_n}{\delta_k}.$$

Similarly, one can prove the inequalities for  $\kappa_k$ .

To prove part 2 we consider the function

$$\psi(t) = \prod_{j=1}^n \frac{t + \gamma_j}{t + \lambda_j} = \prod_{j=1}^n \left( 1 + \frac{\gamma_j - \lambda_j}{t + \lambda_j} \right).$$

Since  $\gamma_j - \lambda_j \geq 0$  for  $j = 1, \dots, n$ , it follows that  $\psi(t)$  is decreasing as  $t$  increases. Since  $\psi(\delta_k) = \epsilon_k$  for  $k = 1, \dots, n$ , and  $\delta_1 < \dots < \delta_n$ , we thus have

$$\epsilon_1 > \epsilon_2 > \dots > \epsilon_n.$$

Obviously  $\psi(t) > 1$  for any  $t > 0$ , and hence  $\epsilon_n = \psi(\delta_n) > 1$ .

The monotonicity  $\kappa_1 > \dots > \kappa_n > 1$  follows in the same way. □

With the help of Lemma 4.1 we can now prove the following entrywise monotonicity property of the minimal positive solution  $X$  of (1.1).

**THEOREM 4.2.** *Let  $X = [x_{ij}] \in \mathbb{R}^{n,n}$  be the minimal positive solution of (1.1). Then for any  $i \geq k$  and  $j \geq l$  with  $(i, j) \neq (k, l)$ , the entries of  $X$  satisfy*

$$x_{ij} > x_{kl}.$$

*Proof.* Since

$$0 < \gamma_1 < \dots < \gamma_n, \quad 0 < \delta_1 < \dots < \delta_n,$$

and by Lemma 4.1,

$$1 < \epsilon_n < \cdots < \epsilon_1, \quad 1 < \kappa_n < \cdots < \kappa_1,$$

with (3.11), for  $1 \leq i, j \leq n$ , if  $i < n$ , it follows that

$$x_{ij} = \frac{\epsilon_i \kappa_j}{\delta_i + \gamma_j} > \frac{\epsilon_{i+1} \kappa_j}{\delta_{i+1} + \gamma_j} = x_{i+1,j}.$$

If  $j < n$ , then

$$x_{ij} = \frac{\epsilon_i \kappa_j}{\delta_i + \gamma_j} > \frac{\epsilon_i \kappa_{j+1}}{\delta_i + \gamma_{j+1}} = x_{i,j+1}. \quad \square$$

The quantities in Lemma 4.1 also provide upper and lower bounds for the entries of the minimal positive solution  $X$  of (1.1).

**THEOREM 4.3.** *Let  $X = [x_{ij}] \in \mathbb{R}^{n,n}$  be the minimal positive solution of (1.1). Then*

$$\frac{w_{ij}}{\delta_i + \gamma_j} < x_{ij} < \frac{W_{ij}}{\delta_i + \gamma_j},$$

where

$$w_{ij} = \max \left\{ \frac{a_i b_j}{p_i p_j}, \frac{a_i}{p_i}, \frac{b_j}{p_j}, 1 \right\},$$

$$W_{ij} = \min \left\{ \frac{\delta_i \gamma_j}{p_i p_j}, \frac{\delta_i (\gamma_j + \delta_n)}{p_i \gamma_j}, \frac{(\delta_i + \gamma_n) \gamma_j}{\delta_i p_j}, \frac{(\delta_i + \gamma_n) (\gamma_j + \delta_n)}{\delta_i \gamma_j} \right\}.$$

*Proof.* The bounds follow from the formulas (3.8)–(3.11) and the inequalities given in the first part of Lemma 4.1.  $\square$

**COROLLARY 4.4.** *Let  $X = [x_{ij}] \in \mathbb{R}^{n,n}$  be the minimal positive solution of (1.1), and let  $w_{ij}, W_{ij}$  be as in Theorem 4.3. Then*

$$\frac{w_{nn}}{\delta_n + \gamma_n} < x_{nn} \leq x_{ij} \leq x_{11} < \frac{W_{11}}{\delta_1 + \gamma_1}$$

for  $i, j = 1, \dots, n$ .

*Proof.* The inequalities follow from Theorems 4.2 and 4.3.  $\square$

By taking advantage of the scaled equation (2.2), we also obtain a bound for the spectral norm of the minimal positive solution  $X$  of (1.1).

**THEOREM 4.5.** *Let  $\tilde{X} \in \mathbb{R}^{n,n}$  be the minimal positive solution of (2.2). Then*

$$\|\tilde{X}\| \leq 1,$$

and  $\|\tilde{X}\| = 1$  if and only if  $\chi(0) = 0$  and  $\chi'(0) = 0$ .

Moreover, the minimal positive solution  $X$  of (1.1) satisfies

$$\|X\| \leq \frac{1}{\min_j p_j}.$$

*Proof.* Let  $\tilde{X}_+ \geq \tilde{X}$  be another positive solution of (2.2) [20]. Since both  $\tilde{X}$  and  $\tilde{X}_+$  are positive, it is easily verified that

$$\|\tilde{X}\|^2 = \rho(\tilde{X}^T \tilde{X}) \leq \rho(\tilde{X}^T \tilde{X}_+),$$

where  $\rho(Z)$  is the spectral radius of  $Z$ . Lemma 3.1 shows that  $\tilde{X}^T$  is the minimal positive solution of the dual equation (3.1). By Lemma 12 of [7],  $\rho(\tilde{X}^T \tilde{X}_+) \leq 1$ . Hence  $\|\tilde{X}\| \leq 1$ , and  $\|\tilde{X}\| = 1$  if and only if  $\tilde{X}_+ = \tilde{X}$ . The last equality holds if and only if 0 is a double eigenvalue of  $\tilde{H}$ , which is equivalent to the conditions  $\chi(0) = 0$  and  $\chi'(0) = 0$ , by Lemma 2.1.

The upper bound for  $\|X\|$  follows from the relation  $X = \Phi^{-1} \tilde{X} \Phi^{-1}$ . □

Various lower bounds for  $\|X\|$  can also be derived by using the inequalities for the entries of  $X$ , but we will not pursue this topic here.

At the end of this section we also provide a formula for the inverse of  $X$ .

**THEOREM 4.6.** *The minimal positive solution  $X = [x_{ij}]$  of (1.1) is invertible and with  $P, \Theta$  as in Theorem 3.3, its inverse is given by*

$$X^{-1} = PQ\Theta^TGP,$$

where

$$Q = \text{diag}(q_1, \dots, q_n), \quad G = \text{diag}(g_1, \dots, g_n),$$

with

$$q_k = \prod_{j=1}^n \frac{\gamma_k + \delta_j}{\gamma_k - \lambda_j}, \quad g_k = \prod_{j=1}^n \frac{\delta_k + \gamma_j}{\delta_k - \nu_j},$$

for  $k = 1, \dots, n$ .

*Proof.* Since  $\gamma_n > \dots > \gamma_1 > 0$  and  $\delta_n > \dots > \delta_1 > 0$ , it follows (see, e.g., [6]) that the Cauchy matrix  $\Theta$  is invertible and

$$\Theta^{-1} = \hat{Q}\Theta^T\hat{G},$$

where

$$\hat{Q} = \text{diag}(\hat{q}_1, \dots, \hat{q}_n), \quad \hat{G} = \text{diag}(\hat{g}_1, \dots, \hat{g}_n),$$

with

$$\hat{q}_k = \frac{\prod_{j=1}^n (\gamma_k + \delta_j)}{\prod_{j \neq k} (\gamma_k - \gamma_j)}, \quad \hat{g}_k = \frac{\prod_{j=1}^n (\delta_k + \gamma_j)}{\prod_{j \neq k} (\delta_k - \delta_j)},$$

for  $k = 1, \dots, n$ . Since all of the diagonal matrices in (3.8) are invertible, it follows that  $X$  is also invertible and the formula for  $X^{-1}$  follows from (3.8) using  $\Theta^{-1}$ . □

**5. Numerical algorithms.** The formulas given in section 3 can be used to develop the following numerical algorithms for computing the minimal positive solution of (1.1).

**ALGORITHM 5.1.** For the Riccati equation (1.1) this algorithm computes the minimal positive solution.

1. Compute the eigenvalues  $-\nu_1, \dots, -\nu_n, \lambda_1, \dots, \lambda_n$  of  $\tilde{H}$  in (2.3) by applying a root finding solver to the secular equation  $\chi(\lambda) = 0$  given by (2.4).
2. Use either of the formulas (3.8) or (3.11) to compute the minimal positive solution  $X$  of (1.1).

We can also use either of the formulas (3.9) or (3.10).

ALGORITHM 5.2. For the Riccati equation (1.1) this algorithm computes the minimal positive solution.

1. Compute the eigenvalues  $-\nu_1, \dots, -\nu_n$  of  $\tilde{H}$  in (2.3) by applying a root finding solver to the secular equation  $\chi(\lambda) = 0$  given by (2.4).
2. Use formula (3.9) to compute the minimal positive solution  $X$  of (1.1).

ALGORITHM 5.3. For the Riccati equation (1.1) this algorithm computes the minimal positive solution.

1. Compute the eigenvalues  $\lambda_1, \dots, \lambda_n$  of  $\tilde{H}$  in (2.3) by applying a secular equation solver to  $\chi(\lambda) = 0$ .
2. Use formula (3.10) to compute the minimal positive solution  $X$  of (1.1).

Note that Algorithms 5.2 and 5.3 need only compute half of the eigenvalues.

The success of these three algorithms depends on how fast and accurately the eigenvalues can be computed and how sensitive the evaluation of the formulas (3.8)–(3.11) is. This requires an efficient and reliable secular equation solver. The osculatory interpolation methods of [3, 24] that were developed in the context of the divide-and-conquer eigenvalue methods ([10, section 8.5], [4, 5, 9]) may not be applicable directly, since the secular function  $\chi(\lambda)$  has quite different properties than the secular equation derived in the symmetric divide-and-conquer method. For this reason we propose the following hybrid method for the computation of roots of the secular function. We consider only the case for computing the eigenvalues  $\lambda_k$  as the method for computing the eigenvalues  $\nu_k$  is analogous. Our approach treats  $\lambda_1$  differently from the other eigenvalues  $\lambda_2, \dots, \lambda_n$ , because of the different properties that  $\lambda_1$  has.

### 5.1. Computation of $\lambda_k$ with $k > 1$ .

1. Initial guess. To compute an initial guess, we basically follow the procedure suggested in [24]. We first evaluate  $\chi(m_k)$ , where  $m_k$  is the midpoint of the interval  $(\gamma_{k-1}, \gamma_k)$ . Because  $\chi(\lambda)$  has only one root in  $(\gamma_{k-1}, \gamma_k)$ , and since  $\lim_{\lambda \rightarrow \gamma_{k-1}^+} \chi(\lambda) = \infty$ , and  $\lim_{\lambda \rightarrow \gamma_k^-} \chi(\lambda) = -\infty$ , based on the sign of  $\chi(m_k)$ , we can easily determine in which half of the interval  $\lambda_k$  is located. Simple geometry shows that if  $\chi(m_k) > 0$ , then  $\lambda_k$  is closer to  $\gamma_k$ , and if  $\chi(m_k) < 0$ , then  $\lambda_k$  is closer to  $\gamma_{k-1}$ . We then consider the equation

$$\frac{p_{k-1}}{\lambda - \gamma_{k-1}} + \frac{p_k}{\lambda - \gamma_k} + r_k = 0,$$

with  $r_k = \chi(m_k) - p_{k-1}/(m_k - \gamma_{k-1}) - p_k/(m_k - \gamma_k)$ , which can be obtained during the evaluation of  $\chi(m_k)$  without any extra cost. We then take the root of this equation in  $(\gamma_{k-1}, \gamma_k)$  as our initial guess  $z_k^0$ . It is easily verified that  $z_k^0$  and  $\lambda_k$  are in the same half interval. We also choose an initial interval so that the  $\chi$  values on endpoints have opposite signs (which guarantees that  $\lambda_k$  is in this interval). If  $\chi(m_k)\chi(z_k^0) < 0$ , then we use  $m_k, z_k^0$  for the interval. Otherwise, we use the asymptotic properties of  $\chi$  to find another  $\lambda$  value to replace  $m_k$ . Let us denote the resulting interval by  $[u_0, v_0]$ .

2. Iteration step. For a current approximation  $z_k^j$ , we first evaluate  $\chi'(z_k^j)$  and use one step of Newton's method to determine the next approximate  $z_k^{j+1}$ . If  $z_k^{j+1}$  is inside the current interval  $[u_j, v_j]$ , then we evaluate  $\chi(z_k^{j+1})$ . We then replace one of  $u_j, v_j$  and its corresponding  $\chi$  value with  $z_k^{j+1}$  and  $\chi(z_k^{j+1})$  based on the sign of  $\chi(z_k^{j+1})$  and move on to the next iteration. If  $z_k^{j+1}$  is outside  $[u_j, v_j]$  (maybe even outside of  $(\gamma_{k-1}, \gamma_k)$ ), then we apply one step of

the secant method with  $u_j, v_j$  and their corresponding  $\chi$  values to get  $z_k^{j+1}$ . We then evaluate  $\chi(z_k^{j+1})$ , update  $[u_j, v_j]$ , and continue. If this  $z_k^{j+1}$  is still outside of  $[u_j, v_j]$ , then we use one step of the bisection method with  $u_j, v_j$  to get  $z_k^{j+1}$ .

When the iterates  $z_k^j$  get close to the root  $\lambda_j$ , then, due to rounding errors, it becomes more difficult to compute a reliable value of  $\chi(z_k^j)$ . (This happens typically for small roots.) This may cause the sign of  $\chi$  to alternate between positive and negative values in the Newton iteration and the secant iteration, which may have the effect that the sequence  $\{z_k^j\}$  does not converge. If we observe such a behavior and the function values for  $\chi$  are also small in absolute value, then we run a step of the bisection method. This procedure has turned out to be very successful during our numerical tests.

3. Stopping criterion. In order to compute the root  $\lambda_k$  accurately, we actually use the shift  $s = \lambda - \gamma_{k-1}$  or  $s = \lambda - \gamma_k$  initially, depending on whether  $\lambda_k$  is closer to  $\gamma_{k-1}$  or  $\gamma_k$ . The iteration step is then applied to the new variable  $s$  to generate a sequence of approximate values  $s_0, s_1, \dots, s_j, \dots$ . The iteration can be written as

$$s_{j+1} = s_j + \Delta s_j,$$

where  $\Delta s_j$  is the  $j$ th correction.

We use the stopping criterion

$$(5.1) \quad |\Delta s_j| < c\varepsilon_M |s_{j+1}|,$$

where  $\varepsilon_M$  is the machine epsilon and  $c$  is a modest constant (which is set to 48 in our tests).

The procedure for the computation of  $\nu_k$  ( $k = 2, \dots, n$ ) is analogous.

### 5.2. Computation of $\lambda_1$ .

1. Initial guess. The strategy for choosing starting values  $z_1^0$  and starting intervals  $[u_0, v_0]$  is slightly different than in the case of the other eigenvalues. Since we know that  $\lambda_1 \in [0, \gamma_1)$ , we first evaluate  $\chi(m_1)$ , where  $m_1 = \gamma_1/2$ . We use the sign of  $\chi(m_1)$  to determine if  $\lambda_1$  is closer to 0 or  $\gamma_1$ . We then use the root  $z_1^0 \in [0, \gamma)$  of the equation

$$\frac{p_1}{\lambda - \gamma_1} + r_1 = 0,$$

with  $r_1 = \chi(m_1) - p_1/(m_1 - \gamma_1)$ , as the initial starting value.

If  $\chi(m_1), \chi(z_1^0) < 0$ , then we use  $m_1, z_1^0$  to form the initial interval  $[u_0, v_0]$ .

If  $\chi(m_1), \chi(z_1^0) > 0$ , then we replace  $m_1$  by another value such that the corresponding  $\chi$  value is negative, by using the fact  $\lim_{\lambda \rightarrow \gamma_1^-} (\lambda) = -\infty$ . In

the case that  $\chi(m_1), \chi(z_1^0) < 0$ , if  $\chi(0) > 0$ , we replace  $m_1$  with 0. If  $\chi(0) = 0$ , we still need to check the sign of  $\chi'(0)$ . If  $\chi'(0) > 0$ , we may use it to find a small positive number such that its corresponding  $\chi$  is positive. We then replace  $m_1$  with this number. If  $\chi'(0) \leq 0$ , we simply set  $\lambda_1 = 0$ , and no iteration is required.

Note that for the transport theory problem,  $\chi(0)$  and  $\chi'(0)$  can be easily determined by the formulas given in Remark 2.2.

2. Iteration step. We first use the same iteration steps as described for the eigenvalues  $\lambda_k$ ,  $k \geq 2$ , to an approximation of  $\lambda_1$ . This usually works well for  $\lambda_1 > c_1 \sqrt{\varepsilon_M}$  with some positive constant  $c_1$ . If, however,  $\lambda_1$  is too small, then it is difficult to get accurate function values for  $\chi$  and  $\chi'$ , which then may cause convergence problems. In order to overcome this difficulty, once we observe that the  $j$ th approximate  $z_1^j$  satisfies  $z_1^j < c_1 \sqrt{\varepsilon_M}$  (we used  $c_1 = 100$  in our tests), we evaluate  $\chi(z_1^j)$  and  $\chi'(z_1^j)$  by using their corresponding Taylor polynomials at 0, given by

$$\chi(z_1^j) \approx \chi(0) + z_1^j \chi'(0) + \frac{(z_1^j)^2}{2} \chi''(0),$$

$$\chi'(z_1^j) \approx \chi'(0) + z_1^j \chi''(0) + \frac{(z_1^j)^2}{2} \chi'''(0),$$

and use these values in the next step of the Newton iteration. If  $\chi'(z_1^j)$  is also very small in modulus, then we approximate  $\chi''(z_1^j)$  by

$$\chi''(z_1^j) \approx \chi''(0) + z_1^j \chi'''(0).$$

We then use the approximations for  $\chi(z_1^j)$ ,  $\chi'(z_1^j)$ ,  $\chi''(z_1^j)$  to construct the second degree Taylor polynomial for  $\chi$  at  $z_1^j$  and use one of the roots of this polynomial (if it exists) as our next iterate  $z_1^{j+1}$ .

For a general secular equation, the computation of  $\chi(0)$ ,  $\chi'(0)$ ,  $\chi''(0)$ , and  $\chi'''(0)$  requires extra cost and it is not clear if the values can be evaluated accurately. In the secular equation from the transport problem, however, this computation is essentially cost-free since we may use the formulas in Remark 2.2, and because of the simple formulations the values can be computed accurately.

3. Stopping criterion. We use again the stopping criterion (5.1) (with  $\gamma_0 := 0$ ). The procedure for the computation of  $\nu_1$  is analogous.

**5.3. Costs.** The main cost in Algorithms 5.1–5.3 is the evaluation of  $\chi$  and  $\chi'$  during each iteration step. In order to evaluate  $\chi(\lambda)$  and  $\chi'(\lambda)$ , we first compute  $\lambda - \gamma_j$ ,  $\lambda + \delta_j$  for  $j = 1, \dots, n$ . We then compute  $p_j/(\lambda - \gamma_j)$  and  $p_j/(\lambda + \delta_j)$ . After this  $\chi(\lambda)$  can be evaluated. We continue to compute  $[p_j/(\lambda - \gamma_j)]/(\lambda - \gamma_j)$  and  $[p_j/(\lambda + \delta_j)]/(\lambda + \delta_j)$ , which costs one extra flop for each term, and then evaluate  $\chi'(\lambda)$ . So if the Newton iteration is used in the iteration step, then the cost per iteration step and per eigenvalue is about  $10n$  flops. If the average number of iterations is  $M$ , then the cost for Algorithm 5.1 is about  $(20M + 9)n^2$  flops, and the cost for Algorithms 5.2 and 5.3 is about  $(10M + 9)n^2$  flops. Note that it requires  $3n^2$  flops to compute each set of the values  $\xi_k, \eta_k, \kappa_k, \epsilon_k$ , and it requires another  $3n^2$  flops to compute the components of  $X$ . Note also that in these complexity estimates we did not count the cost for the computation of the initial values.

**5.4. Error analysis.** To analyze the computational errors in the described procedures, we first estimate the errors in the computed eigenvalues; see also [30]. We assume that the iteration for each eigenvalue stops when (5.1) holds, and the computed sequence satisfies the conditions in the following lemma observed by Kahan (see, e.g., [24]).

LEMMA 5.4. Let  $\{x_j\}_{j=1}^\infty$  be a sequence of real numbers produced by some rapidly convergent iteration scheme, such that  $\lim_{j \rightarrow \infty} x_j = x^*$ . If the sequence of ratios  $\frac{|x_{j+1}-x_j|}{|x_j-x_{j-1}|}$  is decreasing for  $j \geq k$ , and if  $\frac{|x_{k+1}-x_k|}{|x_k-x_{k-1}|} < 1$ , then

$$|x_{k+1} - x^*| < \frac{|x_{k+1} - x_k|^2}{|x_k - x_{k-1}| - |x_{k+1} - x_k|}.$$

Let  $\lambda_j, \nu_j$  be the exact eigenvalues of  $H$ , and let  $\hat{\lambda}_j, \hat{\nu}_j$  be the corresponding computed eigenvalues. With the discussed properties of the eigenvalues, the presented procedures, and Lemma 5.4, it is reasonable to assume that the computed eigenvalues satisfy

$$(5.2) \quad |\lambda_j - \hat{\lambda}_j| < C_{\lambda_j} \varepsilon_M \min\{\gamma_j - \lambda_j, \lambda_j - \gamma_{j-1}\},$$

$$(5.3) \quad |\nu_j - \hat{\nu}_j| < C_{\nu_j} \varepsilon_M \min\{\delta_j - \nu_j, \nu_j - \delta_{j-1}\},$$

for  $j = 1, \dots, n$ , where  $\gamma_0 = \delta_0 = 0$  and  $C_{\lambda_j}, C_{\nu_j}$  are some modest constants. We then have the following lemma.

LEMMA 5.5. Suppose that the computed eigenvalues  $\hat{\lambda}_j, -\hat{\nu}_j$  of  $H$  as in (1.5) satisfy (5.2) and (5.3). Let  $\hat{\xi}_k, \hat{\eta}_k, \hat{\epsilon}_k, \hat{\kappa}_k$  be the computed quantities determined via the formulas given in Theorem 3.3. Then

$$\hat{\xi}_k = \xi_k(1 + nC_{\xi_k} \varepsilon_M), \quad \hat{\eta}_k = \eta_k(1 + nC_{\eta_k} \varepsilon_M),$$

$$\hat{\kappa}_k = \kappa_k(1 + nC_{\kappa_k} \varepsilon_M), \quad \hat{\epsilon}_k = \epsilon_k(1 + nC_{\epsilon_k} \varepsilon_M),$$

for  $k = 1, \dots, n$ , where  $C_{\xi_k}, C_{\eta_k}, C_{\kappa_k}, C_{\epsilon_k}$  are constants.

*Proof.* For the proof we just consider the first order error.

Note that  $\hat{\xi}_k$  is actually computed by the formula

$$\prod_{j=1}^n (\gamma_k - \hat{\lambda}_j) / \prod_{j \neq k} (\gamma_k - \gamma_j);$$

i.e.,  $\lambda_j$  is replaced with  $\hat{\lambda}_j$ . We then have

$$|\gamma_k - \hat{\lambda}_j| = |(\gamma_k - \lambda_j) + (\lambda_j - \hat{\lambda}_j)| = |\gamma_k - \lambda_j| \left| 1 + \frac{\lambda_j - \hat{\lambda}_j}{\gamma_k - \lambda_j} \right| =: |\gamma_k - \lambda_j| |1 + \tilde{C}_{kj} \varepsilon_M|,$$

for  $j = 1, \dots, n$ , where by (5.2) and the interlacing property of the eigenvalues

$$|\tilde{C}_{kj}| = \frac{1}{\varepsilon_M} \left| \frac{\lambda_j - \hat{\lambda}_j}{\gamma_k - \lambda_j} \right| < C_{kj} \frac{\min\{\gamma_j - \lambda_j, \lambda_j - \gamma_{j-1}\}}{|\gamma_k - \lambda_j|} \leq C_{kj}.$$

With this relation, it is not difficult to obtain that

$$\hat{\xi}_k = \xi_k(1 + nC_{\xi_k} \varepsilon_M),$$

where  $C_{\xi_k}$  is a constant. The corresponding relations for the other terms follow in the same way.  $\square$

Using this lemma we obtain the following relative errors for the components of the minimal positive solution computed by the formulas given in section 3.

**THEOREM 5.6.** *Consider the problem of computing the minimal positive solution  $X = [x_{ij}]$  of (1.1) using formulas (3.8)–(3.11), and suppose that the computed eigenvalues satisfy the relations (5.2) and (5.3). Then for the computed solution  $\hat{X} = [\hat{x}_{ij}]$ , the relative error estimate*

$$\frac{|\hat{x}_{ij} - x_{ij}|}{x_{ij}} = D_{ij} n \varepsilon_M, \quad i, j = 1, \dots, n$$

holds, where  $D_{ij}$ 's are positive constants.

*Proof.* The relative error estimates follow from Lemma 5.5.  $\square$

**6. Numerical examples.** In this section we present some numerical test results for the problems from transport theory; see [20, 21]. The weights  $c_1, \dots, c_n$  and nodes  $\omega_1, \dots, \omega_n$  are generated from the composite four-node Gauß–Legendre quadrature formula on  $[0, 1]$  with  $n/4$  equally spaced subintervals; see, e.g., [29]. All the numerical examples were tested in MATLAB version 7.1.0 with machine precision  $\varepsilon_M \approx 2.22e - 16$ . We solved the problem for various numbers of the parameters  $\alpha$  and  $\beta$  and the size  $n$ . We used all four formulas to compute the minimal positive solution, with a secular equation solver as described in section 5.

The computed minimal positive solutions via formulas (3.8)–(3.11) are denoted by  $X^{(1)}, X^{(2)}, X^{(3)}, X^{(4)}$ , respectively. In the following we display the test results. We present one table for each pair  $(\alpha, \beta)$  and various values of  $n$ . (The used norm is always the spectral norm.) In each of Tables 6.1–6.6, we list the following results:

- Maximum residual:

$$R = \max_{j \in \{1, 2, 3, 4\}} \|X^{(j)} \Gamma + \Delta X^{(j)} - (e + X^{(j)} p)(e^T + p^T X^{(j)})\|.$$

- Maximum and minimum entrywise relative errors:

$$RE_{\max} = \max_{\substack{i, j \in \{1, 2, 3, 4\} \\ i \neq j}} \max_{k, l \in \{1, \dots, n\}} \frac{|x_{kl}^{(i)} - x_{kl}^{(j)}|}{\min\{x_{kl}^{(i)}, x_{kl}^{(j)}\}},$$

$$RE_{\min} = \min_{\substack{i, j \in \{1, 2, 3, 4\} \\ i \neq j}} \max_{k, l \in \{1, \dots, n\}} \frac{|x_{kl}^{(i)} - x_{kl}^{(j)}|}{\min\{x_{kl}^{(i)}, x_{kl}^{(j)}\}}.$$

- Largest entry  $x_{11}$  (determined by one of the four solutions).
- Smallest entry  $x_{nn}$  (determined by one of the four solutions).
- Norm  $\|X\|$  ( $X$  is one of the four solutions). Note that we have proved that  $\|\tilde{X}\| \leq 1$ , which translates to  $\|X\| \leq 1/\min p_j$ .
- Number of iterations for  $\nu_1$ :  $N_-$ .
- Number of iterations for  $\lambda_1$ :  $N_+$ .
- Average of the number of iterations for all  $2n$  eigenvalues:  $N$ .

We also give the eigenvalues  $-\nu_1, \lambda_1$  in the caption.

We can summarize the numerical results as follows.

1. The values of  $R$  in the tables are usually the residual of  $X^{(1)}$ . The other residuals are basically the same, but some can be one order smaller.
2. Since we do not know the exact solution, we use  $RE_{\max}$  and  $RE_{\min}$  to detect if high relative accuracy can actually be achieved. The values of  $RE_{\max}$  and  $RE_{\min}$  do support the high relative accuracy result. (Note that  $x_{nn}$  is small in all examples.)



3. The number of iterations for  $\nu_1$  and  $\lambda_1$  increases as  $\alpha \rightarrow 0$  and  $\beta \rightarrow 1$ . This shows the numerical difficulty when the eigenvalues  $-\nu_1$  and  $\lambda_1$  are getting close to each other. However, our computed values of  $\nu_1, \lambda_1$  are much more accurate than those obtained by running the MATLAB code *eig* on  $\tilde{H}$ .
4. Our MATLAB implementation of the root finder based on the secular equation is still not very robust. In general, about .5% of the eigenvalues need 100 iterations, the maximum iteration number used in our experimental code. Some further improvement could enhance these convergence properties.

TABLE 6.1  
 $\alpha = 0.5, \beta = 0.5, (-\nu_1, \lambda_1) \approx (-1.166, 3.996)$ .

$n$	$R$	$RE_{\max}$	$RE_{\min}$	$x_{11}$	$x_{nn}$	$\ X\ $	$N_-$	$N_+$	$N$
64	2.70e-13	1.83e-14	6.80e-15	.263	8.23e-04	7.87e+00	8	7	5
128	1.27e-12	6.72e-14	3.33e-14	.263	4.09e-04	1.57e+01	9	8	5
256	5.35e-12	1.64e-13	7.73e-14	.264	2.04e-04	3.15e+01	9	9	5
512	1.97e-11	2.70e-13	1.34e-13	.264	1.02e-04	6.29e+01	10	8	5

TABLE 6.2  
 $\alpha = 0.1, \beta = 0.99, (-\nu_1, \lambda_1) \approx (-7.98e - 02, 3.83e - 01)$ .

$n$	$R$	$RE_{\max}$	$RE_{\min}$	$x_{11}$	$x_{nn}$	$\ X\ $	$N_-$	$N_+$	$N$
64	5.16e-13	2.65e-14	1.23e-14	2.70	2.19e-03	6.12e+01	8	6	5
128	2.43e-12	9.67e-14	4.06e-14	2.72	1.08e-03	1.22e+02	10	5	5
256	8.48e-12	1.46e-13	7.03e-14	2.72	5.37e-04	2.45e+02	9	5	5
512	3.48e-11	4.21e-13	2.04e-13	2.72	2.67e-04	4.89e+02	10	6	6

TABLE 6.3  
 $\alpha = 10^{-4}, \beta = 1 - 10^{-8}, (-\nu_1, \lambda_1) \approx (-7.91e - 05, 3.79e - 04)$ .

$n$	$R$	$RE_{\max}$	$RE_{\min}$	$x_{11}$	$x_{nn}$	$\ X\ $	$N_-$	$N_+$	$N$
64	2.46e-11	1.48e-12	7.35e-13	4.19	2.24e-03	8.59e+01	23	16	5
128	1.02e-10	5.16e-12	2.57e-12	4.21	1.10e-03	1.72e+02	26	25	5
256	4.66e-11	1.24e-12	5.60e-13	4.22	5.48e-04	3.43e+02	19	25	5
512	5.43e-10	7.02e-12	3.48e-12	4.22	2.73e-04	6.87e+02	34	25	6

TABLE 6.4  
 $\alpha = 10^{-14}, \beta = 1 - 10^{-14}, (-\nu_1, \lambda_1) \approx (-1.73e - 07, 1.73e - 07)$ .

$n$	$R$	$RE_{\max}$	$RE_{\min}$	$x_{11}$	$x_{nn}$	$\ X\ $	$N_-$	$N_+$	$N$
64	6.09e-13	2.52e-14	1.02e-14	4.19	2.24e-03	8.59e+01	28	26	6
128	2.72e-12	7.80e-14	3.15e-14	4.21	1.10e-03	1.72e+02	28	26	5
256	1.02e-11	1.85e-13	8.30e-14	4.22	5.48e-04	3.44e+02	28	26	5
512	4.28e-11	4.12e-13	1.60e-13	4.22	2.73e-04	6.87e+02	28	26	6

TABLE 6.5  
 $\alpha = 10^{-8}, \beta = 1, (-\nu_1, \lambda_1) = (0, 3.00e - 08)$ .

$n$	$R$	$RE_{\max}$	$RE_{\min}$	$x_{11}$	$x_{nn}$	$\ X\ $	$N_-$	$N_+$	$N$
64	7.74e-13	4.84e-14	1.94e-14	4.19	2.24e-03	8.59e+01	0	30	5
128	2.95e-12	8.97e-14	4.07e-14	4.21	1.10e-03	1.72e+02	0	30	5
256	1.21e-11	1.76e-13	7.39e-14	4.22	5.48e-04	3.44e+02	0	32	5
512	4.51e-11	4.14e-13	1.87e-13	4.22	2.73e-04	6.87e+02	0	30	6

TABLE 6.6  
 $\alpha = 10^{-15}$ ,  $\beta = 1$ ,  $(-\nu_1, \lambda_1) = (0, 3.00e - 15)$ .

$n$	$R$	$RE_{\max}$	$RE_{\min}$	$x_{11}$	$x_{nn}$	$\ X\ $	$N_-$	$N_+$	$N$
64	6.97e-13	3.39e-14	1.42e-14	4.19	2.24e-03	8.59e+01	0	55	5
128	2.71e-12	7.83e-14	2.91e-14	4.21	1.10e-03	1.72e+02	0	55	5
256	1.02e-11	1.60e-13	7.47e-14	4.22	5.48e-04	3.44e+02	0	55	5
512	4.19e-11	3.71e-13	1.53e-13	4.22	2.73e-04	6.87e+02	0	55	5

**7. Conclusion.** We have presented four formulas for the minimal positive solution of the nonsymmetric Riccati equation (1.1) that depend on the eigenvalues of the associated matrix. With the help of the formulas we have given some properties and entrywise bounds for the minimal positive solution. We have used the formulas to develop fast numerical algorithms for computing the minimal positive solution. If the eigenvalues can be computed accurately, then the computed minimal positive solution has high relative accuracy.

**Acknowledgments.** We thank an anonymous referee for suggestions that helped to improve this paper. Hongguo Xu wishes to gratefully acknowledge the hospitality of TU Berlin, where part of this research was carried out.

#### REFERENCES

- [1] Z.-Z. BAI, X.-X. GUO, AND S.-F. XU, *Alternately linearized implicit iteration methods for the minimal nonnegative solutions of the nonsymmetric algebraic Riccati equations*, Numer. Linear Algebra Appl., 13 (2006), pp. 655–674.
- [2] D. A. BINI, B. IANNAZZO, AND F. POLONI, *A fast Newton's method for a nonsymmetric algebraic Riccati equation*, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 276–290.
- [3] J. R. BUNCH, C. P. NIELSON, AND D. C. SORESENSEN, *Rank-one modification of the symmetric eigenproblem*, Numer. Math., 31 (1978), pp. 31–48.
- [4] J. J. M. CUPPEN, *A divide and conquer method for the symmetric tridiagonal eigenproblem*, Numer. Math., 36 (1980/81), pp. 177–195.
- [5] J. J. DONGARRA AND D. C. SORESENSEN, *A fully parallel algorithm for the symmetric eigenvalue problem*, SIAM J. Sci. Statist. Comput., 8 (1987), pp. S139–S154.
- [6] T. FINCK, G. HEINIG, AND K. ROST, *An inversion formula and fast algorithms for Cauchy-Vandermonde matrices*, Linear Algebra Appl., 183 (1993), pp. 179–191.
- [7] S. FITAL AND C.-H. GUO, *Convergence of the solution of a nonsymmetric matrix Riccati differential equation to its stable equilibrium solution*, J. Math. Anal. Appl., 318 (2006), pp. 648–657.
- [8] B. D. GANAPOL, *An investigation of a simple transport model*, Transport Theory Statist. Phys., 21 (1992), pp. 1–37.
- [9] G. H. GOLUB, *Some modified matrix eigenvalue problems*, SIAM Rev., 15 (1973), pp. 318–344.
- [10] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, 3rd ed., Johns Hopkins University Press, Baltimore, 1996.
- [11] M. GU AND S. C. EISENSTAT, *A divide-and-conquer algorithm for the symmetric tridiagonal eigenproblem*, SIAM J. Matrix Anal. Appl., 16 (1995), pp. 172–191.
- [12] C.-H. GUO, *Nonsymmetric algebraic Riccati equations and Wiener-Hopf factorization for  $M$ -matrices*, SIAM J. Matrix Anal. Appl., 23 (2001), pp. 225–242.
- [13] C.-H. GUO, *A note on the minimal nonnegative solution of a nonsymmetric algebraic Riccati equation*, Linear Algebra Appl., 357 (2002), pp. 299–302.
- [14] C.-H. GUO AND N. J. HIGHAM, *Iterative solution of a nonsymmetric algebraic Riccati equation*, SIAM J. Matrix Anal. Appl., 29 (2007), pp. 396–412.
- [15] C.-H. GUO, B. IANNAZZO, AND B. MEINI, *On the doubling algorithm for a (shifted) nonsymmetric algebraic Riccati equations*, SIAM J. Matrix Anal. Appl., 29 (2007), pp. 1083–1100.
- [16] C.-H. GUO AND A. J. LAUB, *On the iterative solution of a class of nonsymmetric algebraic Riccati equations*, SIAM J. Matrix Anal. Appl., 22 (2000), pp. 376–391.
- [17] X.-X. GUO, W.-W. LIN, AND S.-F. XU, *A structure-preserving doubling algorithm for nonsymmetric algebraic Riccati equation*, Numer. Math., 103 (2006), pp. 393–412.

- [18] J. JUANG, *Existence of algebraic matrix Riccati equations arising in transport theory*, Linear Algebra Appl., 230 (1995), pp. 89–100.
- [19] J. JUANG AND I.-D. CHEN, *Iterative solution for a certain class of algebraic matrix Riccati equations arising in transport theory*, Transport Theory Statist. Phys., 22 (1993), pp. 65–80.
- [20] J. JUANG AND W.-W. LIN, *Nonsymmetric algebraic Riccati equations and Hamiltonian-like matrices*, SIAM J. Matrix Anal. Appl., 20 (1999), pp. 228–243.
- [21] J. JUANG AND Z. T. LIN, *Convergence of an iterative technique for algebraic matrix Riccati equations and applications to transport theory*, Transport Theory Statist. Phys., 21 (1992), pp. 87–100.
- [22] J. JUANG, C. L. HSING, AND P. NELSON, *Global existence, asymptotics and uniqueness for the reflection kernel of the angularly shifted transport equation*, Math. Models Methods Appl. Sci., 5 (1995), pp. 239–251.
- [23] P. LANCASTER AND L. RODMAN, *Algebraic Riccati Equations*, Oxford University Press, New York, 1995.
- [24] R.-C. LI, *Solving Secular Equations Stably and Efficiently*, Technical Report UCB//CSD-94-851, LAPACK Working Note 93, Computer Science Division, Department of EECS, University of California, Berkeley, CA, 1994.
- [25] L.-Z. LU, *Solution form and simple iteration of a nonsymmetric algebraic Riccati equation arising in transport theory*, SIAM J. Matrix Anal. Appl., 26 (2005), pp. 679–685.
- [26] V. MEHRMANN AND H. XU, *Choosing poles so that the single-input pole placement problem is well conditioned*, SIAM J. Matrix Anal. Appl., 19 (1998), pp. 664–681.
- [27] A. MELMAN, *Numerical solution of a secular equation*, Numer. Math., 69 (1995), pp. 483–493.
- [28] L. C. G. ROGERS, *Fluid models in queueing theory and Wiener-Hopf factorization of Markov Chains*, Ann. Appl. Probab., 4 (1994), pp. 390–413.
- [29] G. W. STEWART, *Afternotes on Numerical Analysis*, SIAM, Philadelphia, 1996.
- [30] G. W. STEWART AND J.-G. SUN, *Matrix Perturbation Theory*, Academic Press, Boston, 1990.