

## Refined rigorous perturbation bounds for the SR decomposition

Mahvish Samar<sup>1,\*</sup>      Aamir Farooq<sup>2</sup>

**Abstract.** In this article, some new rigorous perturbation bounds for the SR decomposition under normwise or componentwise perturbations for a given matrix are derived. Also, the explicit expressions for the mixed and componentwise condition numbers are presented by utilizing the block matrix-vector equation approach. Hypothetical and trial results demonstrate that these new bounds are constantly more tightly than the comparing ones in the literature.

### §1 Introduction

Let  $\mathbb{R}^{m \times n}$  be the arrangement of  $m \times n$  real matrices,  $\mathbb{R}_r^{m \times n}$  be the subset of  $\mathbb{R}^{m \times n}$  comprising of matrices with rank  $r$ ,  $A^T$  be the transpose of matrix  $A$  and  $I_r$  be the identity matrix of order  $r$ . Assume  $A \in \mathbb{R}^{2n \times 2n}$ , and  $P = [e_1, e_3, \dots, e_{2n-1}, e_2, e_4, \dots, e_{2n}]$  with  $e_k$  speaking to the  $k$ -th unit vector. If all even leading principal submatrices of  $PA^TJAP^T$  are nonsingular, Bunse-Gerstner [1] demonstrated that  $A$  has the accompanying SR decomposition

$$A = SR = \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix} \begin{bmatrix} R_{11} & R_{12} \\ R_{21} & R_{22} \end{bmatrix}, \quad (1.1)$$

where  $S \in \mathbb{R}^{2n \times 2n}$  is a symplectic matrix, i.e., it fulfills  $S^TJS = S$ ,

$$J = \begin{bmatrix} 0 & I_n \\ -I_n & 0 \end{bmatrix} \in \mathbb{R}^{2n \times 2n},$$

$R_{ij}(i, j = 1, 2)$  are upper triangular, and  $\text{diag}(R_{21}) = 0$ . Further, if

$$\text{diag}(R_{11}) = |\text{diag}(R_{22})| \quad \text{and} \quad \text{diag}(R_{12}) = 0, \quad (1.2)$$

then the SR decomposition is unique [2]. Throughout this article, we always assume that the factor  $R$  satisfies (1.2).

The SR decomposition is a helpful instrument in the calculation of some optimal control problems (e.g., [3-5]). Also the SR decomposition is a key step for constructing structure-

---

Received: 2020-03-29.      Revised: 2020-11-07.

MR Subject Classification: 65F35, 15A23, 15A57.

Keywords: SR decomposition, Rigorous perturbation bound, Lyapunov majorant function, Banach fixed point theorem, Mixed and componentwise condition numbers.

Digital Object Identifier(DOI): <https://doi.org/10.1007/s11766-021-4086-x>.

The work is supported by the National Natural Science Foundation of China (Grant No. 11771265).

\*Corresponding author.

preserving methods in order to solve the eigenproblem of an important class of structured matrices [2,6-8]. For more points of interest, see for instance [9,10-13]. Henceforth, it is essential to see how perturbations in the original matrix influence consequence of such a decomposition.

The perturbation analysis for the SR decomposition was first introduced by Bhatia [9]. Then Chang [2] also considered this problem and its variants. They both gave first-order perturbation bounds for the SR decomposition. Later, the acquired first-order bounds for SR decomposition was enhanced by Xie et. al [14] and they also presented the rigorous normwise perturbation bounds. However, this bound can severely overestimate the true effect of a perturbation.

In this paper, we investigate some new rigorous perturbation bounds for the SR decomposition under normwise or componentwise perturbations. In addition, the obtained first-order bounds [14] are optimal, which lead to the normwise condition numbers for SR decomposition. However, we know that the normwise condition number may overestimate the ill-posedness of the problem because they ignore the structure of coefficient matrices regarding sparsity or scaling [15,16]. So, it is essential to study the mixed and componentwise condition numbers of SR decomposition. The mixed condition numbers measure the errors in the input data using componentwise error analysis and the output using the normwise error analysis. The componentwise condition numbers measure the errors using the componentwise for both input and output data. Inspired by this, we attempt to present the explicit expressions of the mixed and componentwise condition numbers for the factors  $S$  and  $R$ .

This paper is organized as follows. In Section 3, we combine the block matrix-vector equation approach, the method of Lyapunov majorant function (e.g., [17, Chapter 5]), and the Banach fixed point theorem (e.g., [17 Appendix 5]) to study rigorous perturbation bounds for the SR decomposition when the original matrix has the normwise or componentwise perturbations. In Section 4, we give the explicit expressions for mixed and componentwise condition numbers for the this decomposition. In addition, in Section 2, we provide some notations and preliminaries and in Section 5, we give some numerical experiments.

## §2 Notations and preliminaries

Some notation can be endorsed from [18,19] to make the presentation apparent. We can still illustrate them here to make easier for readers.

For the given matrix  $A = (a_{ij}) \in \mathbb{R}^{m \times n}$ , its spectral norm and Frobenius norm are betoken by  $\|A\|_2$  and  $\|A\|_F$ , respectively. For these two matrix norms, the following inequalities clasp (see [20, pp.80]):

$$\|XYZ\|_2 \leq \|X\|_2 \|Y\|_2 \|Z\|_2, \quad \|XYZ\|_F \leq \|X\|_2 \|Y\|_F \|Z\|_2, \quad (2.1)$$

whenever the matrix product  $XYZ$  is well-defined. For the matrix  $A = (A_{ij}) \in \mathbb{R}^{2n \times 2n}$ , where  $A_{ij} \in \mathbb{R}^{2 \times 2}$ ,  $i, j = 1, 2, \dots, n$ , we define the accompanying operators:

$$\text{upb}(A) = \begin{bmatrix} \frac{1}{2}A_{11} & A_{12} & \cdots & A_{1n} \\ 0 & \frac{1}{2}A_{22} & \cdots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{2}A_{nn} \end{bmatrix}, \quad \text{utb}(A) = \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ 0 & A_{22} & \cdots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & A_{nn} \end{bmatrix},$$

$$\text{uvecb}(A) = \left[ \begin{array}{c} \text{vec}(A_{11}) \\ \vdots \\ \text{vec}(A_{1n}) \\ \hline \text{vec}(A_{22}) \\ \vdots \\ \text{vec}(A_{2n}) \\ \hline \vdots \\ \text{vec}(A_{(n-1)(n-1)}) \\ \text{vec}(A_{(n-1)n}) \\ \hline \text{vec}(A_{nn}) \end{array} \right] \in \mathbb{R}^{v_1}, \quad \text{vecb}(A) = \left[ \begin{array}{c} \text{vec}(A_{11}) \\ \vdots \\ \text{vec}(A_{1n}) \\ \hline \vdots \\ \text{vec}(A_{n1}) \\ \hline \vdots \\ \text{vec}(A_{nn}) \end{array} \right] \in \mathbb{R}^{4n^2},$$

where  $v_1 = 2n(n + 1)$ . Making use of the structures of the operators defined above, we have

$$\begin{aligned} \text{uvecb}(A) &= M_{\text{uvecb}} \text{vecb}(A), & \text{vecb}(\text{utb}(A)) &= M_{\text{utb}} \text{vecb}(A), \\ \text{vecb}(\text{upb}(A)) &= M_{\text{upb}} \text{vecb}(A), \end{aligned} \tag{2.2}$$

where

$$\begin{aligned} M_{\text{uvecb}} &= \text{diag}(S_1, S_2, \dots, S_n) \in \mathbb{R}^{v_1 \times 4n^2}, \quad S_i = [0_{4(n-i+1) \times 4(i-1)}, I_{4(n-i+1)}] \in \mathbb{R}^{4(n-i+1) \times 4n}, \\ M_{\text{utb}} &= \text{diag}(\hat{S}_1, \hat{S}_2, \dots, \hat{S}_n) \in \mathbb{R}^{4n^2 \times 4n^2}, \quad \hat{S}_i = \text{diag}(0_{4(i-1) \times 4(i-1)}, I_{4(n-i+1)}) \in \mathbb{R}^{4n \times 4n}, \\ M_{\text{upb}} &= \text{diag}(\tilde{S}_1, \tilde{S}_2, \dots, \tilde{S}_n) \in \mathbb{R}^{4n^2 \times 4n^2}, \quad \tilde{S}_i = \text{diag}(0_{4(i-1) \times 4(i-1)}, 1/2I_4, I_{4(n-i)}) \in \mathbb{R}^{4n \times 4n}. \end{aligned}$$

Moreover,

$$M_{\text{uvecb}} M_{\text{uvecb}}^T = I_{v_1}, \quad M_{\text{uvecb}}^T M_{\text{uvecb}} = M_{\text{utb}}. \tag{2.3}$$

Thus, letting  $\text{uvecb}^\dagger : \mathbb{R}^{v_1} \rightarrow \mathbb{R}^{2n \times 2n}$  be the right inverse of the operator ‘uvecb’ such that  $\text{uvecb} \cdot \text{uvecb}^\dagger = 1_{v_1 \times v_1}$  and  $\text{uvecb}^\dagger \cdot \text{uvecb} = \text{utb}$ . Then the matrix of the operator ‘uvecb’ is  $M_{\text{uvecb}}^T$ . That is,  $\text{uvecb}^\dagger(A) = M_{\text{uvecb}}^T \text{vecb}(A)$ .

For the operator ‘upb’, the following properties are needed later in this paper. Let  $\mathbb{D}_{2n} \in \mathbb{R}^{2n \times 2n}$  denote the set of diagonal positive definite matrices with  $2 \times 2$  main diagonal blocks  $s_i I_2$ , where  $s_i > 0$ ,  $i = 1, 2, \dots, n$ . Then, for any  $D_{2n} \in \mathbb{D}_{2n}$ ,

$$\text{upb}(AD_{2n}) = \text{upb}(A)D_{2n}, \quad D_{2n}\text{upb}(A) = D_{2n}\text{upb}(A). \tag{2.4}$$

Furthermore, from [2],

$$\|\text{upb}(A) - D_{2n}^{-1}\text{upb}(A^T)D_{2n}\|_F \leq \sqrt{1 + \varsigma_{D_{2n}}^2} \|A\|_F, \tag{2.5}$$

where  $\varsigma_{D_{2n}} = \max_{1 \leq i < j \leq n} \{s_j/s_i\}$ .

To give the definitions of mixed and componentwise condition numbers, the following form of relative distance function will be useful (see [21] for detail). For two vectors  $a = [a_1, \dots, a_p]^T$  and  $b = [b_1, \dots, b_p]^T \in \mathbb{R}^p$ , we define the entry-wise division with

$$c_{i_0} = \begin{cases} \frac{a_{i_0}}{b_{i_0}}, & \text{if } b_{i_0} \neq 0, \\ a_{i_0}, & \text{if } b_{i_0} = 0. \end{cases}$$

Then we define the componentwise distance between  $a$  and  $b$  by

$$d(a, b) = \left\| \frac{a - b}{b} \right\|_\infty = \max_{i=1, \dots, p} \left\{ \frac{|a_{i_0} - b_{i_0}|}{|b_{i_0}|} \right\} = \begin{cases} \frac{|a_{i_0} - b_{i_0}|}{|b_{i_0}|}, & \text{if } b_{i_0} \neq 0, \\ |a_{i_0}|, & \text{if } b_{i_0} = 0. \end{cases}$$

Note that when  $b_{i_0} \neq 0$ ,  $d(a, b)$  will give the relative distance from  $a$  to  $b$  with respect to  $b$ , while the absolute distance for  $b_{i_0} = 0$ . For the distance between the matrices  $A, B \in \mathbb{R}^{n \times n}$ , we define

$$d(A, B) = d(\text{vec}(A), \text{vec}(B)).$$

In addition, for every  $\epsilon > 0$ , we have  $B(a, \epsilon) = \{x \in \mathbb{R}^p \mid |x_i - a_i| \leq \epsilon |a_i|, i = 1, \dots, p\}$ , and denote the domain of definition of function  $F : \mathbb{R}^p \rightarrow \mathbb{R}^q$  as  $\text{Dom}(F)$ .

**Definition 2.1.** [22]: Let  $F : \mathbb{R}^p \rightarrow \mathbb{R}^q$  be a continuous mapping defined on an open set  $\text{Dom}(F) \subset \mathbb{R}^p$ , and  $a \in \text{Dom}(F)$ ,  $a \neq 0$  such that  $F(a) \neq 0$ .

(i) The mixed condition number of  $F$  at  $a$  is defined by

$$m(F, a) = \lim_{\epsilon \rightarrow 0} \sup_{\substack{x \in B(a, \epsilon) \\ x \neq a}} \frac{\|F(x) - F(a)\|_\infty}{\|F(a)\|_\infty} \frac{1}{d(x, a)}.$$

(ii) The componentwise condition number of  $F$  at  $a$  is defined by

$$c(F, a) = \lim_{\epsilon \rightarrow 0} \sup_{\substack{x \in B(a, \epsilon) \\ x \neq a}} \frac{d(F(x), F(a))}{d(x, a)}.$$

When the map  $F$  in Definition 2.1 is Fréchet differentiable, the following lemma given in [22] makes the computation of mixed and componentwise condition number easier.

**Lemma 2.2.** *With the same assumptions as in Definition 2.1, and supposing that  $F$  is Fréchet differentiable at  $a$ , we have*

$$m(F, a) = \frac{\|DF(a)\|_\infty \|a\|_\infty}{\|F(a)\|_\infty}, \quad c(F, a) = \left\| \frac{DF(a)\|a\|}{|F(a)|} \right\|_\infty.$$

where  $DF(a)$  is the Fréchet derivative of  $F$  at  $a$ .

Let  $A = (A_{ij}) \in \mathbb{R}^{2m \times 2n}$  with  $A_{ij} \in \mathbb{R}^{2 \times 2}$ ,  $i = 1, 2, \dots, m$ ,  $j = 1, 2, \dots, n$ . Like the result for the regular operator ‘vec’, the following result holds for ‘vecb’:

$$\hat{\Pi}_{m,n} \text{vecb}(A) = \text{vecb}(A^T), \tag{2.6}$$

where  $\hat{\Pi}_{m,n} = (\Pi_{m,n} \otimes \Pi_{2,2}) \in \mathbb{R}^{4mn \times 4mn}$  with  $\Pi_{m,n} = \sum_{i=1}^m \sum_{j=1}^n (E_{ij} \otimes E_{ij}^T)$ . In these expressions,

$\otimes$  denotes the Kronecker product [23] and the matrix  $E_{ij} \in \mathbb{R}^{m \times n}$  has in the  $(i, j)$ -th position and zeros elsewhere. Given another matrix  $B$ , the block Kronecker product between  $B$  and  $A$  is defined by

$$B \boxtimes A = \begin{bmatrix} B \otimes A_{11} & B \otimes A_{12} & \cdots & B \otimes A_{1n} \\ B \otimes A_{21} & B \otimes A_{22} & \cdots & B \otimes A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ B \otimes A_{m1} & B \otimes A_{m2} & \cdots & B \otimes A_{mn} \end{bmatrix}.$$

For the block Kronecker product, the following results hold [24]

$$\text{vecb}(ACB) = (B^T \boxtimes A) \text{vecb}(C), \tag{2.7}$$

$$\|B \boxtimes A\|_2 = \|B\|_2 \|A\|_2, \tag{2.8}$$

$$(B \boxtimes A)(C \boxtimes G) = (BC \boxtimes AG), \tag{2.9}$$

$$(B \boxtimes A)^{-1} = B^{-1} \boxtimes A^{-1}, \text{ if } B \text{ and } A \text{ are nonsingular.} \tag{2.10}$$

Here, the matrices  $C$  and  $G$  are of suitable orders and are partitioned appropriately.

### §3 Rigorous perturbation bounds

However, it is well known that the first order perturbation bound may lead to erroneous conclusions because they neglect the high order terms. So it is necessary to discuss the rigorous perturbation bound. Next we consider the rigorous perturbation bounds for the factors  $R$  and  $S$  when the original matrix has normwise or componentwise perturbations in the given matrix. First, we give a unique decomposition theorem for a perturbed matrix.

**Theorem 3.1.** *Given  $A \in \mathbb{R}^{2n \times 2n}$ , consider that all even leading submatrices of  $PA^T JAP^T$  are non singular and  $A$  has the unique SR decomposition. If*

$$\|S^T J \Delta A R^{-1}\|_2 < \sqrt{2} - 1, \quad (3.1)$$

then  $A + \Delta A$  has the following unique SR decomposition

$$A + \Delta A = (R + \Delta R)(S + \Delta S). \quad (3.2)$$

**Proof:** Note that  $J^{-1} = J^T = -J$  and  $S$  is nonsingular. Then, left multiplying and right multiplying  $S^T J S = J$  by  $S^T J$  and  $S^{-1} J^T$ , respectively gives  $S J^T S^T = J^T$ , i.e.  $S J S^T = J$ . Now,

$$\begin{aligned} & P(A + \Delta A)^T J(A + \Delta A)P^T \\ &= P(R^T J R + R^T S^T J \Delta A + (\Delta A)^T J S R + (\Delta A)^T J \Delta A)P^T \\ &= P R^T (J + S^T J \Delta A R^{-1} + R^{-T} (\Delta A)^T J S + R^{-T} (\Delta A)^T J \Delta A R^{-1}) R P^T \\ &= P R^T (J + K) R P^T, \end{aligned} \quad (3.3)$$

where  $K = S^T J \Delta A R^{-1} + R^{-T} (\Delta A)^T J S + R^{-T} (\Delta A)^T J \Delta A R^{-1}$ . Taking the spectral norm on  $K$  and using  $J^T J J = J$

$$\begin{aligned} \|K\|_2 &\leq 2\|S^T J \Delta A R^{-1}\|_2 + \|R^{-T} (\Delta A)^T J \Delta A R^{-1}\|_2 \\ &\leq 2\|S^T J \Delta A R^{-1}\|_2 + \|R^{-T} (\Delta A)^T J^T S J S^T J \Delta A R^{-1}\|_2 \\ &\leq 2\|S^T J \Delta A R^{-1}\|_2 + \|S^T J \Delta A R^{-1}\|_2^2 \end{aligned}$$

by (3.1) we have  $\|K\|_2 < 1$ . Using result from [7] we get,  $\|K_{(2i)}\|_2 < 1$  for  $i = 1, 2, \dots, n$ . Furthermore, obviously,  $\|J_{(2i)} K_{(2i)}\|_2 < 1$  also holds. Then  $I_{2i} - J_{(2i)} K_{(2i)}$  is nonsingular, which together with  $(J + K)_{2i} = J_{(2i)} (I_{2i} - J_{(2i)} K_{(2i)})$  shows that  $(J + K)_{(2i)}$  is nonsingular. Noting the structure of  $R$ , we can proof that  $(R^T (J + K) R)_{2i} = R_{2i}^T (J + K)_{2i} R_{2i}$ . Since  $R_{2i}$  is nonsingular, shows that  $(R^T (J + K) R)_{2i}$  is also nonsingular. Thus, observing (3.3), we can get that all even leading principal submatrices of  $P(A + \Delta A)^T J(A + \Delta A)P^T$  are nonsingular. Thus,  $A + \Delta A$  has unique SR decomposition.

**Remark 3.2.** Here we provide the condition under which the perturbed matrix always has the unique SR factorization, while in [2,14] they only simply assume that is true.

#### 3.1. Refined normwise perturbation bound

Assume that the perturbed SR decomposition defined in (3.2) and using the fact that

$$(A + \Delta A)^T J(A + \Delta A) = (R + \Delta R)^T J(R + \Delta R).$$

Using (1.1) and little simplification we will get

$$\begin{aligned} & S^T J \Delta A R^{-1} + R^{-T} (\Delta A)^T J S + R^{-T} (\Delta A)^T J (\Delta A) R^{-1} - R^{-T} (\Delta R)^T J (\Delta R) R^{-1} \\ & = J \Delta R R^{-1} + (J \Delta R R^{-1})^T. \end{aligned} \tag{3.4}$$

In the above equation  $J \Delta R R^{-1}$  is upper triangular, we have

$$\begin{aligned} & J \Delta R R^{-1} \\ & = \text{upb} (S^T J \Delta A R^{-1} + R^{-T} (\Delta A)^T J S) + \text{upb} (R^{-T} (\Delta A)^T J (\Delta A) R^{-1} - R^{-T} (\Delta R)^T J (\Delta R) R^{-1}). \end{aligned} \tag{3.5}$$

Applying the operator ‘vecb’ to (3.5) and (2.7), (2.6) and (2.3) gives

$$\begin{aligned} (R^{-T} \boxtimes J) \text{vecb}(\Delta R) &= M_{\text{upb}} ((R^{-T} \boxtimes S^T J) + (S^T J \boxtimes R^{-T} J) \Pi_{n,n}) \text{vecb}(\Delta A) \\ &+ M_{\text{upb}} (R^{-T} \boxtimes R^{-T}) \text{vecb}((\Delta A)^T J \Delta A - (\Delta R)^T J \Delta R). \end{aligned}$$

As done in [25,26] and (2.10), we can obtain

$$\begin{aligned} \text{vecb}(\Delta R) &= (R^T \boxtimes J^{-1}) M_{\text{upb}} ((R^{-T} \boxtimes S^T J) + (S^T J \boxtimes R^{-T} J) \Pi_{n,n}) \text{vecb}(\Delta A) \\ &+ (R^T \boxtimes J^{-1}) M_{\text{upb}} (R^{-T} \boxtimes R^{-T}) \text{vecb}((\Delta A)^T J \Delta A - (\Delta R)^T J \Delta R). \end{aligned} \tag{3.6}$$

and show that Eq. (3.6) is equivalent to

$$\begin{aligned} \text{uvecb}(\Delta R) &= M_{\text{uvecb}} (R^T \boxtimes J^{-1}) M_{\text{upb}} ((R^{-T} \boxtimes S^T J) + (S^T J \boxtimes R^{-T} J) \Pi_{n,n}) \text{vecb}(\Delta A) \\ &+ M_{\text{uvecb}} (R^T \boxtimes J^{-1}) M_{\text{upb}} (R^{-T} \boxtimes R^{-T}) \text{vecb}((\Delta A)^T J \Delta A - (\Delta R)^T J \Delta R). \end{aligned} \tag{3.7}$$

For simplicity, let

$$G_R = M_{\text{uvecb}} (R^T \boxtimes J^{-1}) M_{\text{upb}} ((R^{-T} \boxtimes S^T J) + (S^T J \boxtimes R^{-T} J) \Pi_{n,n}), \tag{3.8}$$

$$H_R = M_{\text{uvecb}} (R^T \boxtimes J^{-1}) M_{\text{upb}} (R^{-T} \boxtimes R^{-T}). \tag{3.9}$$

Then (3.7) becomes

$$\text{uvecb}(\Delta R) = G_R \text{vecb}(\Delta A) + H_R \text{vecb}((\Delta A)^T J \Delta A - (\Delta R)^T J \Delta R). \tag{3.10}$$

Thus, applying the operator ‘uvecb<sup>†</sup>’ to (3.10) leads

$$\Delta R = \text{uvecb}^\dagger [G_R \text{vecb}(\Delta A) + H_R \text{vecb}((\Delta A)^T J \Delta A - (\Delta R)^T J \Delta R)].$$

The above equation can be written as an operator equation for  $\Delta R$ :

$$\begin{aligned} \Delta R &= \Phi(\Delta R, \Delta A), \\ &= \text{uvecb}^\dagger [G_R \text{vecb}(\Delta A) + H_R \text{vecb}((\Delta A)^T J \Delta A - (\Delta R)^T J \Delta R)]. \end{aligned} \tag{3.11}$$

We will execute the technique of Lyapunov majorant function and the Banach fixed point principle to probe the rigorous perturbation bounds for  $\Delta R$  based on the operator equation (3.11) as done in [19]. To make it easy and clear for readers and for plenum of the method we comprehend the detail process here through some steps which are same as in [19]. Assume that  $Z \in \mathbb{R}^{(2n) \times (2n)}$  is upper triangular with the same structure as that of  $\Delta R$ ,  $\|Z\|_F \leq \rho$  for some  $\rho \geq 0$ , and  $\|\Delta A\|_F = \delta$ . Then it follows from the definition of the operator ‘uvecb<sup>†</sup>’ and (2.1) that

$$\|\Phi(Z, \Delta K)\|_F \leq \|G_R\|_2 \delta + \|H_R\|_2 (\delta^2 + \rho^2). \tag{3.12}$$

From (3.12), we have the Lyapunov majorant function of the operator equation (3.13)

$$h(\rho, \delta) = \|G_R\|_2 \delta + \|H_R\|_2 (\delta^2 + \rho^2),$$

and the Lyapunov majorant equation

$$h(\rho, \delta) = \rho, \text{ i.e., } \|G_R\|_2 \delta + \|H_R\|_2 (\delta^2 + \rho^2) = \rho. \tag{3.13}$$

Assume that  $\delta \in \Omega = \delta \geq 0 : 1 - 4\|H_R\|_2 (\|G_R\|_2 \delta + \|H_R\|_2 \delta^2) \geq 0$ . Then, the Lyapunov majorant equation (3.8) has two nonnegative roots:  $\rho_1(\delta) \leq \rho_2(\delta)$  with

$$\rho_1(\delta) = f(\delta) = \frac{2(\|G_R\|_2 \delta + \|H_R\|_2 \delta^2)}{1 + \sqrt{1 - 4\|H_R\|_2 (\|G_R\|_2 \delta + \|H_R\|_2 \delta^2)}}.$$

Let the set  $B(\delta)$  be

$$B(\delta) = \{Z \in \mathbb{R}^{(2n) \times (2n)} : \text{Having the same structure as that of } \Delta R \text{ and } \|Z\|_F \leq f(\delta)\},$$

which is closed and convex. We can check that the operator  $\Phi(\cdot, \Delta A)$  maps the set  $B(\delta)$  into itself and for  $Z, \tilde{Z} \in B(\delta)$ ,

$$\left\| \Phi(Z, \Delta A) - \Phi(\tilde{Z}, \Delta A) \right\|_F \leq h'_\rho(f(\delta), \delta) \left\| Z - \tilde{Z} \right\|_F.$$

Since the derivative of the function  $h(\rho, \delta)$  relative to  $\rho$  at  $f(\delta)$  satisfies  $h'_\rho(f(\delta), \delta) = 1 - \sqrt{1 - 4\|H_R\|_2 (\|G_R\|_2 \delta + \|H_R\|_2 \delta^2)} < 1$  when  $\delta \in \Omega_1 = \{\delta \geq 0 : 1 - 4\|H_R\|_2 (\|G_R\|_2 \delta + \|H_R\|_2 \delta^2) > 0\}$ . Then the operator  $\Phi(\cdot, \Delta A)$  is contractive on the set  $B(\delta)$  for  $\delta \in \Omega_1$ . Thus, from the Banach fixed point principle, we have that the operator equation (3.11), i.e., the matrix equation (3.4), has a unique solution in the set  $B(\delta)$ . As a result,  $\|\Delta R\|_F \leq f(\delta)$  for  $\delta \in \Omega_1$ . In summary, we have the following main theorem.

**Theorem 3.3.** *With the same assumptions as in Theorem 3.1, if*

$$\|H_R\|_2 (\|G_R\|_2 \delta + \|H_R\|_2 \delta^2) < \frac{1}{4}, \tag{3.14}$$

*then  $A + \Delta A$  has the unique SR factorization (3.2) and*

$$\|\Delta R\|_2 \leq \frac{2(\|G_R\|_2 \delta + \|H_R\|_2 \delta^2)}{1 + \sqrt{1 - 4\|H_R\|_2 (\|G_R\|_2 \delta + \|H_R\|_2 \delta^2)}} \tag{3.15}$$

$$\leq 2(\|G_R\|_2 \delta + \|H_R\|_2 \delta^2) \tag{3.16}$$

$$< (1 + 2\|G_R\|_2) \|\Delta A\|_F. \tag{3.17}$$

**Proof:** It is easy to see that the condition (3.15) is the same as the one in  $\Omega_1$ . Thus, from the discussions before Theorem 3.1, it suffices to obtain the bound (3.17). This can be done by noting (3.14) and the fact

$$2\|H_R\|_2 \|\Delta A\|_F \leq \sqrt{1 + \|G_R\|_2^2} - \|G_R\|_2 < 1,$$

which can be derived from (3.16).

**Remark 3.4.** From (3.15), by omitting the higher order terms, we can get the first order perturbation bound of  $R$  factor

$$\|\Delta R\|_F \leq \|G_R\|_2 \|\Delta A\|_F + O(\|\Delta A\|_F^2). \tag{3.18}$$

In [14], the authors presented the following optimal first order bound for  $R$

$$\|\Delta R\|_F \leq \|(R^T \boxtimes J^T) \mathfrak{D} K_{SR}\|_2 \|\Delta A\|_F + O(\|\Delta A\|_F^2), \tag{3.19}$$

where

$$\mathfrak{D} \equiv \text{diag}(\mathfrak{D}_1, \mathfrak{D}_2, \mathfrak{D}_2, \mathfrak{D}_1) \in \mathbb{R}^{(4n^2) \times (4n^2)}$$

for  $\mathfrak{D}_1$  and  $\mathfrak{D}_2$  see [14] and

$$K_{SR} \equiv (R^{-T} \boxtimes S^T J) - (S^T J \boxtimes R^{-T})\Pi. \tag{3.20}$$

Now we show that the bound (3.18) is the same as (3.19). In fact, according to the above definition and the definition of the operator 'upb', we can check that for any matrix  $X \in \mathbb{R}^{2n \times 2n}$ ,

$$\mathfrak{D}\text{vecb}(X) = \text{vecb}(\text{upb}(X)).$$

Thus, for any matrix  $X \in \mathbb{R}^{2n \times 2n}$ , using (2.7) and (2.2), we have

$$\begin{aligned} & \| (R^T \boxtimes J^T) \mathfrak{D}K_{SR} \text{vecb}(W) \|_2 \\ &= \| (R^T \boxtimes J^T) \mathfrak{D}\text{vecb}(S^T J W R^{-1} - R^{-T} W^T J^T S) \|_2 \\ &= \| (R^T \boxtimes J^T) \text{vecb}(\text{upb}(S^T J W R^{-1} - R^{-T} W^T J^T S)) \|_2 \\ &= \| \text{vecb}(J^T \text{upb}(S^T J W R^{-1} - R^{-T} W^T J^T S) R) \|_2 \\ &= \| (R^T \boxtimes J^T) M_{\text{upb}}(\text{upb}(S^T J W R^{-1} - R^{-T} W^T J^T S)) \|_2 \\ &= \| (R^T \boxtimes J^T) M_{\text{upb}}((R^{-T} \boxtimes S^T J) - (S^T J \boxtimes R^{-T})\Pi) \text{vecb}(W) \|_2. \end{aligned} \tag{3.21}$$

From the definitions of the matrices 'M<sub>utb</sub>' and 'M<sub>upb</sub>', we can verify that

$$M_{\text{utb}}(R^T \boxtimes J^T) M_{\text{upb}} = (R^T \boxtimes J^T) M_{\text{upb}},$$

which together with (3.21) and (2.3) gives

$$\| (R^T \boxtimes J^T) \mathfrak{D}K_{SR} \text{vecb}(W) \|_2 = \| M_{\text{uvecb}}^T G_R \text{vecb}(W) \|_2 = \| G_R \text{vecb}(W) \|_2.$$

Thus, from the definition of spectral norm we get

$$\begin{aligned} \| (R^T \boxtimes J^T) \mathfrak{D}K_{SR} \|_2 &= \max_{\| \text{vecb}(W) \|_2=1} \| (R^T \boxtimes J^T) \mathfrak{D}K_{SR} \text{vecb}(W) \|_2 \\ &= \max_{\| \text{vecb}(W) \|_2=1} \| G_R \text{vecb}(W) \|_2 \\ &= \| G_R \|_2. \end{aligned}$$

So the bounds (3.18) and (3.19) are the same. Therefore, the rigorous bounds in Theorem 3.3 can be regarded as the rigorous versions of the optimal first order perturbation bound given in [14].

**Remark 3.5.** In [14] the following rigorous perturbation bounds were derived by a combination of the classic and refine matrix equation approaches,

$$\| \Delta R \|_F \leq \sqrt{2} \kappa(R^{-1}) \| S \|_2 \| \Delta A \|_F + (3\sqrt{2} + 2\sqrt{3}) \kappa(R^{-1}) \| S \|_2^2 \| \Delta A \|_F^2, \tag{3.22}$$

under the condition

$$\| S^T J \Delta A R^{-1} \|_F < 1/\sqrt{6} + 2.$$

In the above bounds, for a non singular matrix  $Z$ ,  $\kappa(Z)$  denotes its condition number and is defined as  $\kappa(Z) = \| Z \|_2 \| Z^{-1} \|_2$  the set of  $n \times n$  positive diagonal matrices.

Now we will show that our bounds (3.17) is sharper than (3.22)

$$\begin{aligned} \| G_R \|_2 &= \| M_{\text{uvecb}} (R^T \boxtimes J^{-1}) M_{\text{upb}} ((R^{-T} \boxtimes S^T J) + (S^T J \boxtimes R^{-T} J) \Pi_{n,n}) \|_2 \\ &\leq \| R^T \|_2 \| M_{\text{upb}} ((R^{-T} \boxtimes S^T J) + (S^T J \boxtimes R^{-T} J) \Pi_{n,n}) \|_2 \\ &= \| R^T \|_2 \max_{\| \text{vecb}(X) \|_2=1} \| M_{\text{upb}} ((R^{-T} \boxtimes S^T J) + (S^T J \boxtimes R^{-T} J) \Pi_{n,n}) \text{vecb}(X) \|_2 \end{aligned}$$

$$\begin{aligned}
 &= \|R^T\|_2 \max_{\|\text{vecb}(X)\|_2=1} \|M_{\text{upb}} \text{vecb}(S^T JXR^{-1} - (S^T JXR^{-1})^T)\|_2 \\
 &= \|R^T\|_2 \max_{\|\text{vecb}(X)\|_2=1} \|\text{vecb}(\text{upb}(S^T JXR^{-1} - (S^T JXR^{-1})^T))\|_2 \\
 &= \|R^T\|_2 \max_{\|X\|_F=1} \|\text{upb}(S^T JXR^{-1} - (S^T JXR^{-1})^T)\|_F \\
 &\leq \|R^T\|_2 \max_{\|X\|_F=1} \sqrt{2} \|S^T JXR^{-1}\|_F \\
 &\leq \sqrt{2}\kappa(R^{-1}) \|S\|_2.
 \end{aligned}$$

**Remark 3.6.** From (3.2) we have,

$$\Delta A = S\Delta R + \Delta S R + \Delta S \Delta R. \tag{3.23}$$

Postmultiplying by  $R^{-1}$  gives

$$\Delta S = \Delta A R^{-1} - S \Delta R R^{-1} - \Delta S \Delta R R^{-1}. \tag{3.24}$$

Applying the operator ‘vecb’ to (3.24) and using (2.7), (2.6) and (2.3), we have

$$\text{vecb}(\Delta S) \approx (R^{-T} \boxtimes I_{2n}) \text{vecb}(\Delta A) - (R^{-T} \boxtimes S) \text{vecb}(\Delta R) - (R^{-T} \boxtimes I_{2n}) \text{vecb}(\Delta S \Delta R). \tag{3.25}$$

with the help of above equation and (3.6), we have

$$\begin{aligned}
 \text{vecb}(\Delta S) \approx & ((R^{-T} \boxtimes I_{2n}) - (I_{2n} \boxtimes S J^{-1}) M_{\text{upb}} ((R^{-T} \boxtimes S^T J) \\
 & + (S^T J \boxtimes R^{-T} J) \Pi_{n,n})) \text{vecb}(\Delta A) - (R^{-T} \boxtimes I_{2n}) \text{vecb}(\Delta S \Delta R).
 \end{aligned} \tag{3.26}$$

For simplicity, let

$$G_S = ((R^{-T} \boxtimes I_{2n}) - (I_{2n} \boxtimes S J^{-1}) M_{\text{upb}} ((R^{-T} \boxtimes S^T J) + (S^T J \boxtimes R^{-T} J) \Pi_{n,n})). \tag{3.27}$$

Then (3.26) becomes

$$\text{vecb}(\Delta S) \approx G_S \text{vecb}(\Delta A) - (R^{-T} \boxtimes I_{2n}) \text{vecb}(\Delta S \Delta R). \tag{3.28}$$

Facilitate more, utilizing (3.24) and the outcomes on  $R$  consider given Theorem 3.3, we can get the rigorous perturbation bounds for  $S$  factor. The rigorous perturbation bounds for factor  $S$  is bigger than the one given in [14]. So we don’t talk about their detailed derivation. From (3.26), by omitting the higher order terms, we can get the first order perturbation bound of  $S$  factor

$$\|\Delta S\|_F \leq \|G_S\|_2 \|\Delta A\|_F + O(\|\Delta A\|_F^2). \tag{3.29}$$

Similar to the discussions in Remark 4.2, we can confirm that

$$\|(R^{-T} \boxtimes I_{2n}) - (I_{2n} \boxtimes S J^{-1}) \mathfrak{D}K_{SR}\|_2 = \|G_S\|_2. \tag{3.30}$$

So the bound (3.29) is same as the optimal one in [14].

### 3.2. Componentwise perturbation

As done in [27], in the following, we consider the componentwise perturbation in the given matrix with the following perturbation

$$|\Delta A| \leq \epsilon L |A|, \quad L = (l_{ij}) \in \mathbb{R}^{2n \times 2n}, \quad 0 \leq l_{ij} \leq 1, \tag{3.31}$$

where  $\Delta A$  is the perturbation matrix and  $\epsilon \geq 0$  is a small scalar. Now we consider componentwise rigorous perturbation bounds for  $S$  and  $R$  factors by using matrix equation approach. The

derivations of these bounds are similar to the ones in [14].

**Theorem 3.7.** *Given  $A \in \mathbb{R}^{2n \times 2n}$ , consider that all even leading submatrices of  $PA^T J A P^T$  are non singular and  $A$  has the unique SR factorization and  $\Delta A$  satisfies (3.31). If*

$$\| \|S^T \|J \|L \|R^{-1} \| \|_F \text{cond}(R) \epsilon < \frac{1}{\sqrt{6} + 2}, \tag{3.32}$$

then  $A + \Delta A$  has the following unique SR factorization and

$$\begin{aligned} & \| \Delta R \|_F \tag{3.33} \\ & \leq \frac{\inf_{D \in \mathbb{D}_{2n}} \| \|R \|R^{-1} \|D \| \|D^{-1} R \|_2 (\sqrt{2+2\zeta_D^2} \| \|S^T \|J \|L \|R^{-1} \| \|_F + (\sqrt{3} - \sqrt{2}) \| \|S^T \|J \|L \|R^{-1} \| \|_F)}{\sqrt{2} - 1} \epsilon \end{aligned}$$

$$\| \Delta R \|_F \leq (\sqrt{6} + \sqrt{3}) \left( \inf_{D \in \mathbb{D}_{2n}} \| \|R \|R^{-1} \|D \| \|D^{-1} R \|_2 \right) \| \|S^T \|J \|L \|R^{-1} \| \|_F \epsilon \tag{3.34}$$

$$\| \Delta S \|_F \leq (\sqrt{6} + 2\sqrt{3} + 2 + 2\sqrt{2}) \| \|S^T \|J \|L \|S \| \|_F \| \|S \|_2 \text{cond}(R) \epsilon, \tag{3.35}$$

where  $\text{cond}(R) = \| \|R \|R^{-1} \| \|_2$ .

**Proof:** Note that  $J$  is skew-symmetric. Then, (3.5) equation can be rewritten as

$$\begin{aligned} J \Delta R R^{-1} &= \text{upb} (S^T J \Delta A R^{-1} - (S^T J \Delta A R^{-1})^T) \\ &+ \text{upb} (R^{-T} (\Delta A)^T J (\Delta A) R^{-1} - R^{-T} (\Delta R)^T J (\Delta R) R^{-1}). \end{aligned} \tag{3.36}$$

Taking the Frobenius norm on (3.36) and using (2.5) by putting  $D_{2n} = I_{2n}$  yields

$$\begin{aligned} \| \Delta R R^{-1} \|_F &\leq \sqrt{2} \| S^T J \Delta A R^{-1} \|_F + \frac{1}{\sqrt{2}} \| R^{-T} (\Delta A)^T J^T S J S^T J (\Delta A) R^{-1} \|_F \\ &+ \frac{1}{\sqrt{2}} \| R^{-T} (\Delta R)^T J (\Delta R) R^{-1} \|_F \\ \| \Delta R R^{-1} \|_F &\leq \sqrt{2} \| S^T J \Delta A R^{-1} \|_F + \frac{1}{\sqrt{2}} \| S^T J (\Delta A) R^{-1} \|_F^2 + \frac{1}{\sqrt{2}} \| \Delta R R^{-1} \|_F^2. \end{aligned}$$

Therefore,

$$\frac{1}{\sqrt{2}} \| \Delta R R^{-1} \|_F^2 - \| \Delta R R^{-1} \|_F + \sqrt{2} \| S^T J \Delta A R^{-1} \|_F + \frac{1}{\sqrt{2}} \| S^T J (\Delta A) R^{-1} \|_F^2 \geq 0. \tag{3.37}$$

The inequality (3.37) can be considered as a quadratic inequality on  $\| \Delta R R^{-1} \|_F$ . By the assumption (3.32), we have

$$\begin{aligned} \Upsilon &\equiv (-1)^2 - 4 \times \frac{1}{\sqrt{2}} \times \left( \sqrt{2} \| S^T J \Delta A R^{-1} \|_F + \frac{1}{\sqrt{2}} \| S^T J (\Delta A) R^{-1} \|_F^2 \right) \\ &= 1 - 4 \| S^T J \Delta A R^{-1} \|_F - 2 \| S^T J (\Delta A) R^{-1} \|_F^2 \geq 0. \end{aligned}$$

Hence  $\| \Delta R R^{-1} \|_F \leq \frac{1}{\sqrt{2}} (1 - \sqrt{\Upsilon})$  or  $\| \Delta R R^{-1} \|_F \geq \frac{1}{\sqrt{2}} (1 + \sqrt{\Upsilon})$  Since  $\frac{1}{\sqrt{2}} (1 - \sqrt{\Upsilon})$  or  $\frac{1}{\sqrt{2}} (1 + \sqrt{\Upsilon})$  and  $\| \Delta R R^{-1} \|_F$  are all continuous functions of the  $\Delta A$ , and  $\Delta R \rightarrow 0$  as  $\Delta A \rightarrow 0$ , we must have

$$\| \Delta R R^{-1} \|_F \leq \frac{1}{\sqrt{2}} (1 - \sqrt{1 - 4 \| S^T J (\Delta A) R^{-1} \|_F - 2 \| S^T J (\Delta A) R^{-1} \|_F^2}) < \frac{1}{\sqrt{2}}.$$

Postmultiplying (3.36) by  $D$ , where  $D \in \mathbb{D}_{2n}$  and noting  $R = D \bar{R}$

$$\begin{aligned} J \Delta R \bar{R}^{-1} &= \text{upb} (S^T J \Delta A \bar{R}^{-1} - (S^T J \Delta A \bar{R}^{-1})^T) \\ &+ \text{upb} (R^{-T} (\Delta A)^T J (\Delta A) \bar{R}^{-1} - R^{-T} (\Delta R)^T J (\Delta R) \bar{R}^{-1}). \end{aligned}$$

Taking Frobinous norm and using (2.5), we get

$$\begin{aligned} \|\Delta R\bar{R}^{-1}\|_F &\leq \sqrt{1 + \zeta_D^2} \|S^T J \Delta A \bar{R}^{-1}\|_F + \|R^{-T}(\Delta A)^T J^T S J S^T J(\Delta A) \bar{R}^{-1}\|_F \\ &\quad + \|R^{-T}(\Delta R)^T J(\Delta R) \bar{R}^{-1}\|_F \\ &\leq \sqrt{1 + \zeta_D^2} \|S^T J \Delta A \bar{R}^{-1}\|_F + \|S^T J(\Delta A) R^{-1}\|_F \|S^T J(\Delta A) \bar{R}^{-1}\|_F \\ &\quad + \|\Delta R R^{-1}\|_F \|\Delta R \bar{R}^{-1}\|_F, \end{aligned}$$

using  $\|\Delta R R^{-1}\|_F < 1/\sqrt{2}$ , we have

$$\|\Delta R \bar{R}^{-1}\|_F \leq \frac{\sqrt{2 + 2\zeta_D^2} \|S^T J \Delta A \bar{R}^{-1}\|_F + \sqrt{2} \|S^T J(\Delta A) R^{-1}\|_F \|S^T J(\Delta A) \bar{R}^{-1}\|_F}{\sqrt{2} - 1}.$$

With help of this  $\|S^T J(\Delta A) \bar{R}^{-1}\|_F \leq \|R\| \|\bar{R}^{-1}\|_2 \|S^T\| \|J\| \|L\| \|R^{-1}\|_F \epsilon$ ,  $\sqrt{1 + \zeta_D^2} \geq 1$ , which further reduces to

$$\begin{aligned} &\|\Delta R \bar{R}^{-1}\|_F \\ &\leq \frac{\sqrt{2 + 2\zeta_D^2} \|R\| \|\bar{R}^{-1}\|_2 \|S^T\| \|J\| \|L\| \|R^{-1}\|_F \epsilon + (\sqrt{3} - \sqrt{2}) \|R\| \|\bar{R}^{-1}\|_2 \|S^T\| \|J\| \|L\| \|R^{-1}\|_F \epsilon}{\sqrt{2} - 1} \\ &\leq (\sqrt{3} + \sqrt{6}) \sqrt{1 + \zeta_D^2} \|R\| \|\bar{R}^{-1}\|_2 \|S^T\| \|J\| \|L\| \|R^{-1}\|_F \epsilon. \end{aligned}$$

Noting the fact that

$$\|\Delta R\|_F = \|\Delta R \bar{R}^{-1} \bar{R}\|_F \leq \|\Delta R \bar{R}^{-1}\|_F \|\bar{R}^{-1}\|_2,$$

we will get (3.34). Next we prove (3.35). For this, premultiplying (3.24) by  $S^T J$  and noting  $S^T J S = J$  leads to

$$S^T J \Delta S = S^T J \Delta A R^{-1} - J \Delta R R^{-1} - S^T J \Delta S \Delta R R^{-1}.$$

With the help of (3.36)

$$\begin{aligned} S^T J \Delta S &= S^T J \Delta A R^{-1} - \text{upb}(S^T J \Delta A R^{-1} - R^{-T}(\Delta A)^T J S) \\ &\quad - \text{upb}(R^{-T}(\Delta A)^T J(\Delta A) R^{-1} + R^{-T}(\Delta R)^T J(\Delta R) R^{-1}) - S^T J \Delta S \Delta R R^{-1}. \end{aligned}$$

The above equation can rewritten as

$$\begin{aligned} S^T J \Delta S &= \text{lowb}(S^T J \Delta A R^{-1}) + (\text{lowb}(S^T J \Delta A R^{-1}))^T \\ &\quad - \text{upb}(R^{-T}(\Delta A)^T J^T S J S^T(\Delta A) R^{-1} + R^{-T}(\Delta R)^T J(\Delta R) R^{-1}) \\ &\quad - S^T J \Delta S \Delta R R^{-1}. \end{aligned}$$

Taking the Frobenius norm and using (2.5) yields

$$\begin{aligned} \|S^T J \Delta S\|_F &\leq \sqrt{2} \|S^T J \Delta A R^{-1}\|_F + \frac{1}{\sqrt{2}} \|S^T J \Delta A R^{-1}\|_F^2 + \frac{1}{\sqrt{2}} \|\Delta R R^{-1}\|_F^2 \\ &\quad + \|S^T J \Delta S\|_F \|\Delta R R^{-1}\|_F. \end{aligned} \tag{3.38}$$

Now,

$$\|\Delta R R^{-1}\|_F \leq \frac{1}{\sqrt{2}} (4 \|S^T J(\Delta A) R^{-1}\|_F + 2 \|S^T J(\Delta A) R^{-1}\|_F^2),$$

which combined with  $\|S^T J(\Delta A) R^{-1}\|_F \leq \sqrt{3/2} - 1$  gives

$$\|\Delta R R^{-1}\|_F^2 \leq (5 + 2\sqrt{6}) \|S^T J(\Delta A) R^{-1}\|_F^2. \tag{3.39}$$

Substituting (3.39) into (3.38) and using  $\|\Delta RR^{-1}\|_F < 1/\sqrt{2}$ , we have

$$\begin{aligned} \|S^T J \Delta S\|_F &\leq \frac{\sqrt{2}\|S^T J \Delta A R^{-1}\|_F + (3\sqrt{2} + 2\sqrt{3})\|S^T J \Delta A R^{-1}\|_F^2}{1 - \|\Delta RR^{-1}\|_F} \\ \|S^T J \Delta S\|_F &\leq \frac{(\sqrt{2} + (3\sqrt{2} + 2\sqrt{3}))\|S^T J \Delta A R^{-1}\|_F \|S^T J \Delta A R^{-1}\|_F}{1 - \frac{1}{\sqrt{2}}} \\ \|S^T J \Delta S\|_F &\leq \frac{2 + \sqrt{6}}{\sqrt{2} - 1} \|S^T J \Delta A R^{-1}\|_F. \end{aligned} \tag{3.40}$$

Substituting (3.39), (3.40) in (3.38) we have

$$\begin{aligned} \|S^T J \Delta S\|_F &\leq \sqrt{2}\|S^T J \Delta A R^{-1}\|_F + (3\sqrt{2} + 2\sqrt{3})\|S^T J \Delta A R^{-1}\|_F^2 \\ &\quad + \frac{(2 + \sqrt{6})(\sqrt{3} + \sqrt{2})}{\sqrt{2} - 1} \|S^T J \Delta A R^{-1}\|_F^2. \end{aligned} \tag{3.41}$$

By (3.31), we have

$$\|S^T J \Delta A R^{-1}\|_F \leq \|S^T J \Delta A R^{-1}\|_F \leq \|S^T\| \|J\| \|L\| \|S\| \|R\| \|R^{-1}\|_F \epsilon, \tag{3.42}$$

Now using,  $\|S^T J(\Delta A)R^{-1}\|_F \leq \sqrt{3/2} - 1$  and (3.42) the (3.41) becomes

$$\|S^T J \Delta S\|_F \leq (2\sqrt{2} + 2\sqrt{3} + 2 + \sqrt{6}) \|S^T\| \|J\| \|L\| \|S\| \|R\| \|R^{-1}\|_F \epsilon.$$

Since  $\|\Delta S\|_F = \|S J^T S^T J \Delta S\|_F \leq \|S\|_2 \|S^T J \Delta S\|_F$ . So, we will get (3.35).

Utilizing the methodology of block matrix equation, we can likewise acquire the rigorous perturbation bounds for  $S$  and  $R$  factors with componentwise perturbation i.e. (3.31) however we can't demonstrate that the acquired bounds are always tighter than the (3.33), (3.34) and (3.35). Along these lines, these outcomes are not displayed in this article.

### §4 Mixed and componentwise condition numbers

In this section, we present the explicit expressions for the mixed and componentwise condition numbers. To obtain the explicit expressions for mixed and componentwise condition numbers of SR decomposition (1.1), we need to define the following mappings:

$$\Psi_R : \text{vecb}(A) \rightarrow \text{uvecb}(R), \tag{4.1}$$

$$\Psi_S : \text{vecb}(A) \rightarrow \text{vecb}(S). \tag{4.2}$$

and use the following two first-order approximations:

$$\text{uvecb}(\Delta R) = G_R(\Delta A) + O(\|\Delta A\|_F^2), \tag{4.3}$$

$$\text{vecb}(\Delta S) = G_S(\Delta A) + O(\|\Delta A\|_F^2). \tag{4.4}$$

which are from (3.10) and (3.28), respectively

**Lemma 4.1.** *Given  $A \in \mathbb{R}^{2n \times 2n}$ , consider that all even leading submatrices of  $PA^T J A P^T$  are non singular and  $A$  has the unique SR decomposition. The mapping defined by (4.1) and (4.2) are Fréchet differentiable, and the derivative of  $\Psi_R$  and  $\Psi_S$  at  $a = (\text{vecb}(A)^T)^T$  is given by*

$$D\Psi_R(a) = G_R, \tag{4.5}$$

$$D\Psi_S(a) = G_S. \tag{4.6}$$

where  $G_R$  and  $G_S$  are given in (3.8) and (3.27).

**Proof:** Utilizing the meaning in (4.3), (4.4) and by the definition of Fréchet derivative, we

have the desired results. The definitions for the mixed and componentwise condition numbers for the factors  $R$  and  $S$  are first given as follows

$$\begin{aligned}
 m_R(\Psi_R, a) &= \lim_{\epsilon \rightarrow 0} \sup_{|\Delta A| \leq \epsilon |A|} \frac{\|\Delta R\|_{\max}}{d(a + \Delta a, a) \|R\|_{\max}} = m_R, \\
 c_R(\Psi_R, a) &= \lim_{\epsilon \rightarrow 0} \sup_{|\Delta A| \leq \epsilon |A|} \frac{1}{d(a + \Delta a, a)} \|\Delta R/R\|_{\max} = c_R, \\
 m_S(\Psi_S, a) &= \lim_{\epsilon \rightarrow 0} \sup_{|\Delta A| \leq \epsilon |A|} \frac{\|\Delta S\|_{\max}}{d(a + \Delta a, a) \|S\|_{\max}} = m_S, \\
 c_S(\Psi_S, a) &= \lim_{\epsilon \rightarrow 0} \sup_{|\Delta A| \leq \epsilon |A|} \frac{1}{d(a + \Delta a, a)} \|\Delta S/S\|_{\max} = c_S.
 \end{aligned}$$

where  $a = (\text{vecb}(A)^T)^T$ , and for a matrix  $A$ ,  $\|A\|_{\max} = \|\text{vecb}(A)\|_{\infty} = \max_{i,j} |a_{ij}|$ . Thus, considering Lemma 2.2 and Lemma 4.1, we have the accompanying theorem which gives the explicit expressions of mixed and componentwise condition numbers for the factors  $R$  and  $S$ .

**Theorem 4.2.** *Given  $A \in \mathbb{R}^{2n \times 2n}$ , consider that all even leading submatrices of  $PA^T J A P^T$  are non singular and  $A$  has the unique SR decomposition. Then mixed and componentwise condition numbers for the factors  $R$  and  $S$  are given by*

$$m_R = \frac{\|G_R \text{vecb}(|A|)\|_{\infty}}{\|R\|_{\max}}, \quad c_R = \left\| \frac{|G_R \text{vecb}(|A)|}{|R|} \right\|_{\infty}, \tag{4.7}$$

$$m_S = \frac{\|G_S \text{vecb}(|A|)\|_{\infty}}{\|S\|_{\max}}, \quad c_S = \left\| \frac{|G_S \text{vecb}(|A)|}{|S|} \right\|_{\infty}. \tag{4.8}$$

where  $G_R$  and  $G_S$  are given in (3.8) and (3.27).

In accompanying, we give the upper bounds for the above two condition numbers because the matrices in the exact expressions are very large and will be expensive to compute them.

**Corollary 4.3.** *With the same assumptions as in Theorem 4.2, we have*

$$m_R \leq \frac{\|\text{upb}((|R^{-T}| |A^T| |J| |S|) + (|S^T| |J| |A| |R^{-1}|)) \|R\|_{\max}}{\|R\|_{\max}} = m_R^{upp}(A), \tag{4.9}$$

$$c_R \leq \|\text{upb}((|R^{-T}| |A^T| |J| |S|) + (|S^T| |J| |A| |R^{-1}|))\|_{\max} = c_R^{upp}(A), \tag{4.10}$$

$$m_S \leq \frac{\| |R^{-1}| |A| + |S| |J^{-1}| \text{upb}((|R^{-T}| |A^T| |J| |S|) + (|S^T| |J| |A| |R^{-1}|)) \|_{\max}}{\|S\|_{\max}} = m_S^{upp}(A), \tag{4.11}$$

$$c_S \leq \left\| \frac{|R^{-1}| |A| + |S| |J^{-1}| \text{upb}((|R^{-T}| |A^T| |J| |S|) + (|S^T| |J| |A| |R^{-1}|))}{|S|} \right\|_{\max} = c_S^{upp}(A). \tag{4.12}$$

**Proof:** As for (4.9), it can be obtained from (4.7) and

$$\begin{aligned}
 & \| |M_{\text{uvecb}}(R^T \boxtimes J^{-1}) M_{\text{upb}}((R^{-T} \boxtimes S^T J) + (S^T J \boxtimes R^{-T} J) \Pi_{n,n}) \text{vecb}(|A|) \|_{\infty} \\
 & \leq \| M_{\text{uvecb}}(|R^T| \boxtimes |J^{-1}|) M_{\text{upb}}((|R^{-T}| \boxtimes |S^T J|) + (|S^T J| \boxtimes |R^{-T} J|) \Pi_{n,n}) \text{vecb}(|A|) \|_{\infty} \\
 & = \| \text{uvecb}(|J^{-1}| \text{upb}((|R^{-T}| |A^T| |J| |S|) + (|S^T| |J| |A| |R^{-1}|)) \|R\|) \|_{\infty} \\
 & \leq \| \text{upb}((|R^{-T}| |A^T| |J| |S|) + (|S^T| |J| |A| |R^{-1}|)) \|R\|_{\max}.
 \end{aligned}$$

Similarly, we can obtain (4.10), (4.11) and (4.12).



Table 2. Comparison of rigorous normwise perturbation bounds (3.17) and (3.22).

$n$		$b_{(3.17)}$	$t_{(3.17)}$	$b_{(3.22)}$	$t_{(3.22)}$
2	Pascal	40.6367	0.111850	174.5033	0.015743
	Hilbert	1.2804e+03	0.135869	4.9493e+04	0.023751
	Frank	41.4635	0.132276	155.6911	0.016750
	Random	28.3318	0.085248	646.7755	0.030598
3	Pascal	601.7476	0.159316	4.8353e+03	0.025492
	Hilbert	1.0225e+05	0.140354	1.2674e+08	0.024408
	Frank	1.3537e+03	0.164466	5.2178e+04	0.115959
	Random	2.1857e+03	0.148285	4.9745e+04	0.025030
4	Pascal	3.1207e+07	0.163765	8.5424e+08	0.025447
	Hilbert	4.2752e+08	0.146814	1.7700e+12	0.025031
	Frank	2.7770e+05	0.153419	4.1753e+06	0.026466
	Random	4.8270e+03	0.229636	1.9533e+05	0.024839
5	Pascal	4.0686e+11	0.167354	8.2880e+13	0.028043
	Hilbert	9.4404e+11	0.215623	4.0918e+16	0.025678
	Frank	1.5548e+10	0.282169	9.8352e+12	0.027103
	Random	6.8858e+06	0.346150	6.7399e+08	0.025810
6	Pascal	2.7992e+14	0.318038	9.9907e+16	0.075091
	Hilbert	1.3927e+19	0.218312	1.3368e+24	0.028824
	Frank	3.4876e+13	0.393546	8.0500e+16	0.032341
	Random	7.8562e+09	0.437335	3.3139e+12	0.026591
8	Pascal	1.4795e+20	0.214910	6.0466e+22	0.075685
	Hilbert	2.9537e+32	0.250255	3.4004e+37	0.045839
	Frank	1.1544e+30	0.547057	3.5986e+35	0.037279
	Random	4.1896e+12	0.672344	4.1544e+15	0.042988
10	Pascal	9.7643e+20	0.452109	3.5216e+23	0.085685
	Hilbert	1.7852e+33	0.576208	7.0921e+38	0.056583
	Frank	8.9764e+30	0.862061	1.8352e+36	0.041279
	Random	6.8043e+13	0.663286	9.0526e+16	0.062988

Table 3. Rigorous componentwise perturbation bounds (3.34) and (3.35).

$l$	$n$	4	6	8	10	12	14	30	50
-1	$b_{(3.34)}$	8.2312	12.3658	20.3615	23.4827	30.2538	31.8950	43.6521	76.4218
	$b_{(3.35)}$	39.2562	55.5677	136.9707	141.922	203.3562	294.5619	487.2063	673.65431
-2	$b_{(3.34)}$	0.7990	1.1089	1.8772	2.3197	2.9432	3.3237	22.6706	37.4325
	$b_{(3.35)}$	3.6962	5.0394	12.6822	13.5271	20.1710	21.3739	55.9432	87.3452
-3	$b_{(3.34)}$	0.0822	0.1132	0.1886	0.2219	0.2929	0.3249	18.4312	29.7439
	$b_{(3.35)}$	0.3898	0.5063	1.2157	1.3213	2.0074	2.9681	42.5643	65.0945
-4	$b_{(3.34)}$	0.0066	0.0110	0.0186	0.0215	0.0290	0.0321	12.9812	15.6205
	$b_{(3.35)}$	0.0316	0.0503	0.1248	0.1260	0.1963	0.1989	19.5219	25.0645
-5	$b_{(3.34)}$	0.0005	0.0012	0.0185	0.0022	0.0028	0.0034	5.5409	9.3482
	$b_{(3.35)}$	0.0038	0.0050	0.1155	0.1270	0.1280	0.1310	11.8543	17.4762

From Table 3, we can see that the bound (3.34) is constantly more tightly at that point the bound (3.35).

### References

[1] A Bunse-Gerstner. *Matrix factorization for symplectic QR-like methods*, Linear Algebra Appl, 1986, 83: 49-77.

- [2] X W Chang. *On the sensitivity of the SR decomposition*, Linear Algebra Appl, 1998, 282(1-3): 297-310.
- [3] A Bunse-Gerstner, V Mehrmann. *A symplectic QR-like algorithm for the solution of the real algebraic Riccati equation*, IEEE Trans Automat Control, 1986, 31(12): 1104-1113.
- [4] H A Kandil, G Freiling, V Ionescu, G Jank. *Matrix Riccati Equations in Control and Systems Theory*, Birkhauser Verlag, Basel, 2003.
- [5] V Mehrmann. *The Autonomous Linear Quadratic Control Problem*, Springer-Verlag, Berlin, 1991.
- [6] C Paige, C Van Loan. *A Schur decomposition for Hamiltonian matrices*, Linear Algebra Appl, 1981, 41: 11-32.
- [7] C Van Loan. *A symplectic method for approximating all the eigenvalues of a Hamiltonian matrix*, Linear Algebra Appl, 1984, 61: 233-251.
- [8] D S Watkins, L Elsner. *Self-similar flows*, Linear Algebra Appl, 1983, 110: 213-242.
- [9] R Bhatia. *Matrix factorizations and their perturbations*, Linear Algebra Appl, 1994, 197-198: 245-276.
- [10] A Salam. *On theoretical and numerical aspects of symplectic Gram-Schmidt-like algorithms*, Numer. Algorithms, 2005, 39: 437-462.
- [11] A Salam, E Al-Aidarous. *Error analysis and computational aspects of SR factorization via optimal symplectic householder transformations*, Electron Trans Numer Anal, 2009, 33: 189-206.
- [12] A Salam, E Al-Aidarous, A El Farouk. *Optimal symplectic Householder transformations for SR decomposition*, Linear Algebra Appl, 2008, 429: 1334-1353.
- [13] A Salam, A El Farouk, E Al-Aidarous. *Symplectic Householder transformations for a QR-like decomposition, a geometric and algebraic approaches*, J Comput Appl Math, 2008, 214: 533-548.
- [14] Z Xie, W Li. *Sensitivity analysis for the SR decomposition*, Linear and Multilinear Algebra, 2015, 63: 222-234.
- [15] N J Higham. *Accuracy and Stability of Numerical Algorithms*, second ed, SIAM, Philadelphia, 2002.
- [16] J Sun. *Condition numbers of algebraic Riccati equations in the Frobenius norm*, Linear Algebra Appl, 2002, 350: 237-261.
- [17] M Konstantinov, D Gu, V Mehrmann, P Petkov. *Perturbation Theory for Matrix Equations*, Elsevier, Amsterdam, 2003.
- [18] X W Chang. *Perturbation analysis of some matrix factorization*, Ph D diss, McGill University, 1997.
- [19] Li Hanyu, Y Wei, Y Yang. *New rigorous perturbation bounds for the Cholesky-like factorization of skew-symmetric matrix*, Linear Algebra Appl, 2016, 491: 83-100.
- [20] G W Stewart, J Sun. *Matrix Perturbation Theory*, Academic Press, Boston, 1990.
- [21] Z Xie, W Li, X Jin. *On condition numbers for the canonical generalized polar decomposition of real matrices*, Electron J Linear Algebra, 2013, 26: 842-857.
- [22] F Cucker, H Diao, Y Wei. *On mixed and componentwise condition numbers for Moore-Penrose inverse and linear least squares problems*, Math Comp, 2007, 76: 947-963.

- [23] R A Horn, C R Johnson. *Topics in Matrix Analysis*, Cambridge University Press, Cambridge, 1991.
- [24] R H Koning, H Neudecker, T Wansbeek. *Block Kronecker products and the vecb operator*, Linear Algebra Appl, 1991, 149: 165-184.
- [25] M Samar, A Farooq, H Li, C Mu. *New rigorous perturbation bounds for the generalized Cholesky factorization*, Appl Math Comput, 2019, 362: 124556.
- [26] M Samar, A Farooq, C Mu. *Structured condition numbers and statistical condition estimation for the LDU factorization*, Appl Math J Chin Univ, 2020, 35: 332-348.
- [27] S Singer, S Singer. *Rounding error and perturbation bounds for the symplectic QR factorization*, Linear Algebra Appl, 2003, 358: 255-279.

<sup>1</sup>Department of Mathematics, Shantou University, Shantou 515063, China.

Email: mahvishsamar@hotmail.com

<sup>2</sup>College of Mathematics and Statistics, Chongqing University, Chongqing 401331, China.