



ELSEVIER

Available online at www.sciencedirect.com

 ScienceDirect

Linear Algebra and its Applications 428 (2008) 469–478

LINEAR ALGEBRA
AND ITS
APPLICATIONS

www.elsevier.com/locate/laa

A note on the scaled total least squares problem [☆]

Wei Xu ^a, Sanzheng Qiao ^{a,1}, Yimin Wei ^{b,*}

^a *Department of Computing and Software, McMaster University, Hamilton, Ont., Canada L8S 4K1*

^b *School of Mathematical Sciences, Fudan University, Shanghai 200433, PR China*

Received 25 September 2006; accepted 8 March 2007

Available online 6 April 2007

Submitted by M.K. Ng

Abstract

In this note, we present two results on the scaled total least squares problem. First, we discuss the relation between the scaled total least squares and the least squares problems. We derive an upper bound for the difference between the scaled total least squares solution and the least squares solution and establish a quantitative relation between the scaled total least squares residual and the least squares residual. Second, we give a perturbation analysis of the scaled total least squares problem. Numerical experiments in comparing our results with existing results are demonstrated.

© 2007 Elsevier Inc. All rights reserved.

AMS classification: 15A18; 65F20; 65F25; 65F50

Keywords: Scaled total least squares; Total least squares; Least squares; Perturbation analysis

1. Introduction

The scaled total least squares (STLS) problem is a generalization of the total least squares (TLS) problem. For given $A \in \mathbb{R}^{m \times n}$ ($m > n$) and $b \in \mathbb{R}^m$, the TLS problem is to find $E \in \mathbb{R}^{m \times n}$ and $r \in \mathbb{R}^m$ solving the problem

[☆] The first and second authors are partially supported by the Natural Sciences and Engineering Research Council of Canada. The third author is supported by the National Natural Science Foundation of China under grant 10471027 and Shanghai Education Committee.

* Corresponding author.

E-mail addresses: xuw5@mcmaster.ca (W. Xu), qiao@mcmaster.ca (S. Qiao), ymwei@fudan.edu.cn (Y. Wei).

¹ This work is partially supported by the Shanghai Key Laboratory of Contemporary Applied Mathematics of Fudan University during Sanzheng Qiao's visit.

$$\min_{(b-r) \in \text{range}(A+E)} \|[E \ r]\|_F. \tag{1.1}$$

The STLS generalizes the TLS by introducing a scaling factor. Rao [6] proposed the STLS problem as

$$\min_{(b-r) \in \text{range}(A+E)} \|[E \ \lambda r]\|_F \quad \text{for } E \in \mathbb{R}^{m \times n} \text{ and } r \in \mathbb{R}^m,$$

where $\lambda > 0$ is a given scalar. Obviously, the TLS is a special case of the STLS when $\lambda = 1$. Alternatively, Paige and Strakoš [5] suggested the formulation:

$$\min_{(\lambda b-r) \in \text{range}(A+E)} \|[E \ r]\|_F. \tag{1.2}$$

If $[E_{\text{STLS}} \ r_{\text{STLS}}]$ solves the above problem, then the solution x_{STLS} for x in the equation $(A + E_{\text{STLS}})\lambda x = \lambda b - r_{\text{STLS}}$ is called the STLS solution.

As described above, the relation between the STLS and the TLS is obvious. What is the relation between the STLS and the least squares (LS) problem $\min \|b - Ax\|_2$?

In this note, after giving explicit expressions for the STLS solution x_{STLS} in Section 2, we derive an upper bound for $\|x_{\text{STLS}} - x_{\text{LS}}\|_2$, where x_{LS} denotes the LS solution, and establish a relation between the residuals $\bar{r}_{\text{STLS}} = b - Ax_{\text{STLS}}$ and $r_{\text{LS}} = b - Ax_{\text{LS}}$ in Section 3. Then in Section 4 we present a perturbation analysis of the STLS problem. Finally, in Section 5, we demonstrate our numerical experiments in comparing our bounds with existing ones.

2. Solving STLS problem

In this section, we give existence conditions and explicit expressions for the STLS solution. From the formulation (1.2), if x_{STLS} is the solution of (1.2), then λx_{STLS} is the solution of the TLS problem with A and λb .

The following theorem by Wei [10] gives existence conditions and explicit expressions for the TLS solution.

Theorem 2.1 (Theorem 2.2, [10]). *Let*

$$\check{C} := [A \ b] = \check{U} \check{\Sigma} \check{V}^T \tag{2.1}$$

be the SVD of $[A \ b]$, where $\check{\Sigma} = \text{diag}(\sigma_1(\check{C}), \dots, \sigma_{n+1}(\check{C}))$, $\sigma_1(\check{C}) \geq \dots \geq \sigma_{n+1}(\check{C}) \geq 0$ and $\check{U} \in \mathbb{R}^{m \times (n+1)}$ and $\check{V} \in \mathbb{R}^{(n+1) \times (n+1)}$ have orthonormal columns. Let $k = \text{rank}(A)$, then $\text{rank}(\check{C}) = k + 1$ assuming $b \notin \text{range}(A)$. Partitioning $\check{\Sigma}$, \check{U} , and \check{V} in (2.1):

$$\check{\Sigma} = \begin{bmatrix} \check{\Sigma}_1 & 0 \\ 0 & \check{\Sigma}_2 \end{bmatrix} \begin{matrix} k \\ n - k + 1 \end{matrix}, \quad \check{U} = \begin{bmatrix} \check{U}_1 & \check{U}_2 \\ \check{U}_3 & \check{U}_4 \end{bmatrix} \begin{matrix} k \\ n - k + 1 \end{matrix} \tag{2.2}$$

and

$$\check{V} = \begin{bmatrix} \check{V}_{11} & \check{V}_{12} \\ \check{v}_{21}^T & \check{v}_{22}^T \end{bmatrix} \begin{matrix} n \\ 1 \end{matrix} \tag{2.3}$$

If the conditions

- (i) $\sigma_k(\check{C}) > \sigma_{k+1}(\check{C}) > \sigma_{k+2}(\check{C}) = \dots = \sigma_{n+1}(\check{C}) = 0$,
- (ii) \check{V}_{11} is of full column rank, or equivalently, $\check{v}_{22} \neq 0$,

are satisfied, then

$$\begin{aligned} x_{\text{TLS}} &= (\check{V}_{11}^T)^\dagger \check{v}_{21} = \check{V}_{11} \check{v}_{21} / (1 - \check{v}_{21}^T \check{v}_{21}) \\ &= -\check{V}_{12} (\check{v}_{22}^T)^\dagger = -\check{V}_{12} \check{v}_{22} / (1 - \check{v}_{21}^T \check{v}_{21}) \\ &= (A^T A - \check{V}_{12} \check{\Sigma}_2^2 \check{V}_{12}^T)^\dagger (A^T b - \check{V}_{12} \check{\Sigma}_2^2 \check{v}_{22}) \end{aligned}$$

is the minimal norm TLS solution. Moreover, let $q = \check{v}_{22} / \|\check{v}_{22}\|_2$, then

$$[E_{\text{TLS}} \ r_{\text{TLS}}] = \check{U}_2 \check{\Sigma}_2 q q^T [\check{V}_{12}^T \ \check{v}_{22}]$$

solves (1.1) and

$$\|[E_{\text{TLS}} \ r_{\text{TLS}}]\|_F = \sigma_{k+1}(\check{C}).$$

For the STLS problem, following the formulation (1.2), we consider the SVD

$$C := [A \ \lambda b] = U \Sigma V^T, \tag{2.4}$$

where $U \in \mathbb{R}^{m \times (n+1)}$ has orthonormal columns, $V \in \mathbb{R}^{(n+1) \times (n+1)}$ is orthogonal, and $\Sigma = \text{diag}(\sigma_1(C), \dots, \sigma_{n+1}(C))$, $\sigma_1(C) \geq \dots \geq \sigma_{n+1}(C) \geq 0$. Applying Theorem 2.1, substituting b in (1.1) with λb , and partitioning Σ , U , and V in the SVD (2.4) of C as $\check{\Sigma}$, \check{U} , and \check{V} in (2.2) and (2.3), we can express the STLS solution as

$$\lambda x_{\text{STLS}} = (V_{11}^T)^\dagger v_{21} = -V_{12} (v_{22}^T)^\dagger = (A^T A - V_{12} \Sigma_2^2 V_{12}^T)^\dagger (\lambda A^T b - V_{12} \Sigma_2^2 v_{22}), \tag{2.5}$$

provided that $\sigma_k(C) > \sigma_{k+1}(C)$ and $v_{22} \neq 0$.

Thus, the STLS problem can be solved by the SVD using, for example, $\lambda x_{\text{STLS}} = -V_{12} (v_{22}^T)^\dagger$. In this case, since only V_{12} and v_{22} in V are required, a complete SVD is unnecessary. The SVD can be replaced by any of its approximations as long as a good approximation of the last $n - k + 1$ columns of the V matrix can be obtained. For example, the complete orthogonal decomposition (COD) [2] can be used in place of the SVD. In 1993, Van Huffel and Zha proposed a rank revealing ULV decomposition (RRULVD) method [3]. Although such method is more efficient than the SVD method, its accuracy depends on the estimator for the smallest singular value and its corresponding singular vector. In Section 4, we use the RRULVD method for our perturbation analysis.

3. Relating STLS to LS

While the TLS is a special case of the STLS when $\lambda = 1$, the relation between the STLS and LS is not so obvious. In [4], it is shown that $x_{\text{STLS}} = x_{\text{LS}}$ and $\sigma_{k+1}(C)/\lambda = \|r_{\text{LS}}\|_2$ as $\lambda \rightarrow 0$. In this section, we present quantitative comparisons between the solutions and residuals of the STLS and the LS. Specifically, we derive upper bounds for $\|x_{\text{STLS}} - x_{\text{LS}}\|_2$ and $\|\bar{r}_{\text{STLS}}\|_2$ in terms of $\|r_{\text{LS}}\|_2$.

Theorem 3.1. *If the existence conditions (i) and (ii) in Theorem 2.1 are satisfied, then*

$$\begin{aligned} \|x_{\text{STLS}} - x_{\text{LS}}\|_2 &\leq \rho^2 \|V_{12}^T x_{\text{STLS}} - \lambda^{-1} v_{22}\|_2 + \beta \|x_{\text{STLS}}\|_2 \\ &\leq \frac{\rho^2 + \beta}{\lambda \|v_{22}\|_2}, \end{aligned}$$

where

$$\rho = \frac{\sigma_{k+1}(C)}{\sigma_k(A)} \quad \text{and} \quad \beta = \min \left(1, \frac{\rho^2}{1 - \rho^2} \right). \tag{3.1}$$

Also, the residual norm

$$\|\bar{r}_{\text{STLS}}\|_2 \leq \|r_{\text{LS}}\|_2 + \frac{\rho^2 \sigma_k(A)}{\lambda \|v_{22}\|_2}.$$

Proof. First, we show some equalities used in our derivation. Using the partitions of Σ , U , and V in the SVD (2.4) of C , we can verify

$$A^T A = V_{11} \Sigma_1^2 V_{11}^T + V_{12} \Sigma_2^2 V_{12}^T, \quad \lambda A^T b = V_{11} \Sigma_1^2 v_{21} + V_{12} \Sigma_2^2 v_{22}, \tag{3.2}$$

and

$$V_{12}^T V_{12} + v_{22} v_{22}^T = I. \tag{3.3}$$

From the generalized inverse theory [8], we have

$$(A^T A)^\dagger A^T = A^\dagger \quad \text{and} \quad (I - A^\dagger A) A^T = 0. \tag{3.4}$$

Then, using the first equation in (3.4), $x_{\text{LS}} = A^\dagger b = (A^T A)^\dagger A^T b$, and the second equation in (3.4), we get

$$\begin{aligned} & x_{\text{STLS}} - x_{\text{LS}} \\ &= (I - A^\dagger A) x_{\text{STLS}} + (A^T A)^\dagger V_{12} \Sigma_2^2 V_{12}^T x_{\text{STLS}} + (A^T A)^\dagger (A^T A) x_{\text{STLS}} \\ &\quad - (A^T A)^\dagger V_{12} \Sigma_2^2 V_{12}^T x_{\text{STLS}} - (A^T A)^\dagger A^T b \\ &= (A^T A)^\dagger [(A^T A - V_{12} \Sigma_2^2 V_{12}^T) x_{\text{STLS}} - \lambda^{-1} V_{11} \Sigma_1^2 v_{21}] - \lambda^{-1} (A^T A)^\dagger V_{12} \Sigma_2^2 v_{22} \\ &\quad + (I - A^\dagger A) x_{\text{STLS}} + (A^T A)^\dagger V_{12} \Sigma_2^2 V_{12}^T x_{\text{STLS}}. \end{aligned}$$

From the first equation in (3.2) and $\lambda x_{\text{STLS}} = (V_{11}^T)^\dagger v_{21}$ in (2.5), the expression in the square bracket in the above equation:

$$\begin{aligned} & (A^T A - V_{12} \Sigma_2^2 V_{12}^T) x_{\text{STLS}} - \lambda^{-1} V_{11} \Sigma_1^2 v_{21} \\ &= \lambda^{-1} V_{11} \Sigma_1^2 V_{11}^T (V_{11}^T)^\dagger v_{21} - \lambda^{-1} V_{11} \Sigma_1^2 v_{21} \\ &= 0, \end{aligned}$$

since, by Theorem 2.1, V_{11} is of full column rank and $V_{11}^T (V_{11}^T)^\dagger = I$. Thus

$$x_{\text{STLS}} - x_{\text{LS}} = (I - A^\dagger A) x_{\text{STLS}} + (A^T A)^\dagger V_{12} \Sigma_2^2 (V_{12}^T x_{\text{STLS}} - \lambda^{-1} v_{22}). \tag{3.5}$$

In the following, we show that the first term in the right side of (3.5) satisfies $\|(I - A^\dagger A) x_{\text{STLS}}\|_2 \leq \beta \|x_{\text{STLS}}\|_2$, where β is defined in (3.1).

On the one hand, $\|(I - A^\dagger A) x_{\text{STLS}}\|_2 \leq \|x_{\text{STLS}}\|_2$ since $I - A^\dagger A$ is an orthogonal projection. On the other hand, (2.5) and the symmetry of $A^T A - V_{12} \Sigma_2^2 V_{12}^T$ imply that

$$\begin{aligned} x_{\text{STLS}} &= (A^T A - V_{12} \Sigma_2^2 V_{12}^T)^\dagger (A^T A - V_{12} \Sigma_2^2 V_{12}^T) x_{\text{STLS}} \\ &= (A^T A - V_{12} \Sigma_2^2 V_{12}^T) (A^T A - V_{12} \Sigma_2^2 V_{12}^T)^\dagger x_{\text{STLS}}. \end{aligned}$$

Hence, from the second equation in (3.4),

$$\begin{aligned} & \|(I - A^\dagger A)x_{\text{STLS}}\|_2 \\ &= \|(I - A^\dagger A)(A^T A - V_{12}\Sigma_2^2 V_{12}^T)(A^T A - V_{12}\Sigma_2^2 V_{12}^T)^\dagger x_{\text{STLS}}\|_2 \\ &= \|(I - A^\dagger A)V_{12}\Sigma_2^2 V_{12}^T(A^T A - V_{12}\Sigma_2^2 V_{12}^T)^\dagger x_{\text{STLS}}\|_2 \\ &\leq \|V_{12}\Sigma_2^2 V_{12}^T\|_2 \|(A^T A - V_{12}\Sigma_2^2 V_{12}^T)^\dagger\|_2 \|x_{\text{STLS}}\|_2 \\ &\leq \sigma_{k+1}^2(C) \|(A^T A - V_{12}\Sigma_2^2 V_{12}^T)^\dagger\|_2 \|x_{\text{STLS}}\|_2. \end{aligned}$$

Now, we claim that

$$\|(A^T A - V_{12}\Sigma_2^2 V_{12}^T)^\dagger\|_2 \leq \frac{1}{\sigma_k^2(A) - \sigma_{k+1}^2(C)},$$

then we have $\|(I - A^\dagger A)x_{\text{STLS}}\|_2 \leq \beta \|x_{\text{STLS}}\|_2$. Indeed, from the first equation in (3.2), $A^T A - V_{12}\Sigma_2^2 V_{12}^T$ is of rank k , so

$$\|(A^T A - V_{12}\Sigma_2^2 V_{12}^T)^\dagger\|_2 = \frac{1}{\sigma_k(A^T A - V_{12}\Sigma_2^2 V_{12}^T)}.$$

From Mirsky theorem [7, p. 204], we have

$$\sigma_k(A^T A - V_{12}\Sigma_2^2 V_{12}^T) - \sigma_k(A^T A) \geq -\|V_{12}\Sigma_2^2 V_{12}^T\|_2 \geq -\sigma_{k+1}^2(C)$$

and consequently

$$\|(A^T A - V_{12}\Sigma_2^2 V_{12}^T)^\dagger\|_2 = \frac{1}{\sigma_k(A^T A - V_{12}\Sigma_2^2 V_{12}^T)} \leq \frac{1}{\sigma_k^2(A) - \sigma_{k+1}^2(C)}.$$

For the second term in the right side of (3.5), from $\lambda x_{\text{STLS}} = -V_{12}(v_{22}^T)^\dagger$ in (2.5) and (3.3), we have

$$\begin{aligned} & V_{12}^T x_{\text{STLS}} - \lambda^{-1} v_{22} \\ &= -\lambda^{-1} (V_{12}^T V_{12} (v_{22}^T)^\dagger + v_{22}) \\ &= -\lambda^{-1} (V_{12}^T V_{12} + v_{22} v_{22}^T) v_{22} / (v_{22}^T v_{22}) \\ &= -\lambda^{-1} (v_{22}^T)^\dagger, \end{aligned}$$

which implies

$$\begin{aligned} & \|(A^T A)^\dagger V_{12}\Sigma_2^2 (V_{12}^T x_{\text{STLS}} - \lambda^{-1} v_{22})\|_2 \\ &\leq \rho^2 \|V_{12}^T x_{\text{STLS}} - \lambda^{-1} v_{22}\|_2 \\ &= \frac{\rho^2}{\lambda} \|v_{22}^\dagger\|_2 = \frac{\rho^2}{\lambda \|v_{22}\|_2}. \end{aligned} \tag{3.6}$$

Putting things together, we get

$$\|x_{\text{STLS}} - x_{\text{LS}}\|_2 \leq \frac{\rho^2}{\lambda \|v_{22}\|_2} + \beta \|x_{\text{STLS}}\|_2 \leq \frac{\rho^2 + \beta}{\lambda \|v_{22}\|_2},$$

since $\|\lambda x_{\text{STLS}}\|_2 = \|V_{12}(v_{22}^T)^\dagger\|_2 \leq \|v_{22}^\dagger\|_2 = \|v_{22}\|_2^{-1}$.

Finally, using (3.5) and (3.6), we get the residual norm

$$\begin{aligned} & \|b - Ax_{\text{STLS}}\|_2 \\ & \leq \|b - Ax_{\text{LS}}\|_2 + \|A(A^T A)^\dagger V_{12} \Sigma_2^2 (V_{12}^T x_{\text{STLS}} - \lambda^{-1} v_{22})\|_2 \\ & \leq \|r_{\text{LS}}\|_2 + \frac{\sigma_{k+1}^2(C)}{\lambda \sigma_k(A)} \|v_{22}^\dagger\|_2 = \|r_{\text{LS}}\|_2 + \frac{\rho^2 \sigma_k(A)}{\lambda \|v_{22}\|_2}. \end{aligned}$$

This completes the proof. \square

Since $\sigma_{k+1}(C)/\lambda = \|r_{\text{LS}}\|_2$ as $\lambda \rightarrow 0$ [4], we have $(\rho^2 + \beta)/\lambda \rightarrow 0$ as $\lambda \rightarrow 0$. The above theorem then implies that $x_{\text{STLS}} = x_{\text{LS}}$ and $\|r_{\text{STLS}}\|_2 = \|r_{\text{LS}}\|_2$ when $\lambda \rightarrow 0$, which is consistent with the results in [4].

4. Perturbation analysis

Consider the RRULVD method for solving the STLS problem [3]. Let

$$C := [A \ \lambda b] = P_C \begin{bmatrix} L_C & 0 \\ H_C & F_C \end{bmatrix} Q_C^T,$$

be an RRULVD of C , where L_C is lower triangular and of order $k + 1$ and H_C and F_C are small blocks introduced by rounding errors and approximations. This can also be viewed as a perturbed COD of C . In the RRULVD method, we actually compute the STLS solution of the truncated decomposition:

$$P_C \begin{bmatrix} L_C & 0 \\ 0 & 0 \end{bmatrix} Q_C^T =: [\hat{A} \ \lambda \hat{b}] =: \hat{C}.$$

In this section, we derive an upper bound for $\|x_{\text{STLS}} - \hat{x}_{\text{STLS}}\|_2$, where x_{STLS} and \hat{x}_{STLS} are the STLS solutions of C and \hat{C} respectively. Our analysis can be readily applied to the SVD method, where L_C is diagonal, $H_C = 0$, and F_C a small diagonal matrix.

Since H_C and F_C are introduced by rounding errors, we assume that

$$\Delta C := C - \hat{C} = -P_C \begin{bmatrix} 0 & 0 \\ H_C & F_C \end{bmatrix} Q_C^T$$

is small, specifically,

$$\|H_C\|_2 + \|F_C\|_2 = c \mathbf{u} \|C\|_2 =: \eta, \tag{4.1}$$

where c is a moderate constant and \mathbf{u} is the unit of roundoff.

Before deriving the error bound, it is necessary to verify the existence condition (i) in Theorem 2.1. From (4.1), it follows that

$$\begin{aligned} & \sigma_k(\hat{A}) - \sigma_{k+1}(\hat{C}) \\ & = \sigma_k(A) - \sigma_{k+1}(C) + \sigma_k(\hat{A}) - \sigma_k(A) + \sigma_{k+1}(C) - \sigma_{k+1}(\hat{C}) \\ & \geq \sigma_k(A) - \sigma_{k+1}(C) - 2\eta. \end{aligned}$$

Thus, if $\sigma_k(A) - \sigma_{k+1}(C) > 2\eta$, then the existence condition $\sigma_k(\hat{A}) > \sigma_{k+1}(\hat{C})$ for the perturbed STLS problem is satisfied.

Now, we derive the error bound. Using the SVD (2.4) of C and the partitions of U , Σ , and V , we define

$$E_A := A - U_2 \Sigma_2 V_{12}^T = U_1 \Sigma_1 V_{11}^T \quad \text{and} \quad \lambda e_b := \lambda b - U_2 \Sigma_2 v_{22} = U_1 \Sigma_1 v_{21}.$$

Then, from (2.5), it can be verified that

$$\lambda x_{\text{STLS}} = (V_{11}^T)^\dagger v_{21} = \lambda E_A^\dagger e_b. \tag{4.2}$$

Note that when $\sigma_k(A) > \sigma_{k+1}(C)$, V_{11} is of full column rank [10], implying that $I = V_{11}^\dagger V_{11} = V_{11}^T (V_{11}^T)^\dagger$. Consequently,

$$E_A x_{\text{STLS}} = U_1 \Sigma_1 V_{11}^T x_{\text{STLS}} = \lambda^{-1} U_1 \Sigma_1 V_{11}^T (V_{11}^T)^\dagger v_{21} = \lambda^{-1} U_1 \Sigma_1 v_{21} = e_b.$$

Similarly, letting $\widehat{C} = \widehat{U} \widehat{\Sigma} \widehat{V}^T$ be the SVD of \widehat{C} , partitioning \widehat{U} , $\widehat{\Sigma}$, and \widehat{V} according to (2.2) and (2.3), and defining

$$E_{\widehat{A}} := \widehat{A} - \widehat{U}_2 \widehat{\Sigma}_2 \widehat{V}_{12}^T = \widehat{U}_1 \widehat{\Sigma}_1 \widehat{V}_{11}^T \quad \text{and} \quad \lambda e_{\widehat{b}} := \lambda \widehat{b} - \widehat{U}_2 \widehat{\Sigma}_2 \widehat{v}_{22} = \widehat{U}_1 \widehat{\Sigma}_1 \widehat{v}_{21},$$

we have the solution

$$\widehat{x}_{\text{STLS}} = E_{\widehat{A}}^\dagger e_{\widehat{b}}. \tag{4.3}$$

Comparing the two solutions (4.2) and (4.3), we get

$$\begin{aligned} \widehat{U}_1^T E_{\widehat{A}} (x_{\text{STLS}} - \widehat{x}_{\text{STLS}}) &= \widehat{U}_1^T (E_{\widehat{A}} x_{\text{STLS}} - E_{\widehat{A}} \widehat{x}_{\text{STLS}} + E_A x_{\text{STLS}} - E_A x_{\text{STLS}}) \\ &= \widehat{U}_1^T (E_{\widehat{A}} - E_A) x_{\text{STLS}} - \widehat{U}_1^T (e_{\widehat{b}} - e_b). \end{aligned}$$

Taking the norm on the both sides, we obtain

$$\begin{aligned} \sigma_k(E_{\widehat{A}}) \|x_{\text{STLS}} - \widehat{x}_{\text{STLS}}\|_2 &\leq \| \widehat{U}_1^T E_{\widehat{A}} (x_{\text{STLS}} - \widehat{x}_{\text{STLS}}) \|_2 \\ &\leq \| \widehat{U}_1^T (E_{\widehat{A}} - E_A) \|_2 \|x_{\text{STLS}}\|_2 + \| \widehat{U}_1^T (e_{\widehat{b}} - e_b) \|_2. \end{aligned} \tag{4.4}$$

Obviously, from $E_{\widehat{A}} = \widehat{A} - \widehat{U}_2 \widehat{\Sigma}_2 \widehat{V}_{12}^T$, we get

$$\sigma_k(E_{\widehat{A}}) \geq \sigma_k(\widehat{A}) - \| \widehat{U}_2 \widehat{\Sigma}_2 \widehat{V}_{12}^T \|_2 \geq \sigma_k(\widehat{A}) - \sigma_{k+1}(\widehat{C}) \geq \sigma_k(A) - \sigma_{k+1}(C) - 2\eta, \tag{4.5}$$

since $\text{rank}(E_{\widehat{A}}) = k$. Furthermore, we have

$$\begin{aligned} \| \widehat{U}_1^T (E_{\widehat{A}} - E_A) \|_2 &= \| \widehat{U}_1^T (\widehat{A} - A - \widehat{U}_2 \widehat{\Sigma}_2 \widehat{V}_{12}^T + U_2 \Sigma_2 V_{12}^T) \|_2 \\ &\leq \| \widehat{A} - A \|_2 + \| \widehat{U}_1^T (\widehat{U}_2 \widehat{\Sigma}_2 \widehat{V}_{12}^T + U_2 \Sigma_2 V_{12}^T) \|_2 \\ &\leq \| \widehat{C} - C \|_2 + \| \widehat{U}_1^T U_2 \Sigma_2 V_{12}^T \|_2 \\ &\leq \eta + \sigma_{k+1}(C) \frac{\sigma_{k+1}(\widehat{C}) \eta}{\sigma_{k+1}^2(\widehat{C}) - \eta^2} \end{aligned} \tag{4.6}$$

and, similarly,

$$\begin{aligned} \| \widehat{U}_1^T (e_{\widehat{b}} - e_b) \|_2 &= \| \widehat{U}_1^T (\widehat{b} - b - \lambda^{-1} \widehat{U}_2 \widehat{\Sigma}_2 \widehat{v}_{22} + \lambda^{-1} U_2 \Sigma_2 v_{22}) \|_2 \\ &\leq \| \widehat{b} - b \|_2 + \lambda^{-1} \| \widehat{U}_1^T U_2 \Sigma_2 v_{22} \|_2 \\ &\leq \eta + \lambda^{-1} \sigma_{k+1}(C) \frac{\sigma_{k+1}(\widehat{C}) \eta}{\sigma_{k+1}^2(\widehat{C}) - \eta^2}, \end{aligned} \tag{4.7}$$

since $\| \widehat{U}_1^T U_2 \|_2 \leq \sigma_{k+1}(\widehat{C}) \|H_C\|_2 / (\sigma_{k+1}^2(\widehat{C}) - \|F_C\|_2^2)$ [1, Corollary 2.5] and $\|H_C\|_2, \|F_C\|_2 \leq \eta$. Applying the above three inequalities (4.5), (4.6), and (4.7), from (4.4), we obtain

$$\begin{aligned} & \|x_{\text{STLS}} - \hat{x}_{\text{STLS}}\|_2 \\ & \leq \frac{1}{\sigma_k(E_{\hat{A}})} (\|\widehat{U}_1^T (E_{\hat{A}} - E_A)\|_2 \|x_{\text{STLS}}\|_2 + \|\widehat{U}_1^T (e_{\hat{b}} - e_b)\|_2) \\ & \leq \frac{\eta((1 + \|x_{\text{STLS}}\|_2)(\sigma_{k+1}^2(\widehat{C}) - \eta^2) + \sigma_{k+1}(\widehat{C})\sigma_{k+1}(C)(\|x_{\text{STLS}}\|_2 + \lambda^{-1}))}{(\sigma_k(A) - \sigma_{k+1}(C) - 2\eta)(\sigma_{k+1}^2(\widehat{C}) - \eta^2)}. \end{aligned}$$

The above argument is valid for any sufficiently small perturbation ΔC . Ignoring η^2 , we have the following theorem.

Theorem 4.1. *Let $C = [A \ \lambda b]$. Suppose that $\widehat{C} = C + \Delta C =: [\widehat{A} \ \lambda \widehat{b}]$ and $\|\Delta C\|_2 \approx c\mathbf{u}\|C\|_2 =: \eta \ll 1$, where c is a moderate constant and \mathbf{u} is the unit of roundoff. Let x_{STLS} and \hat{x}_{STLS} be the STLS solutions corresponding to C and \widehat{C} respectively, then*

$$\|x_{\text{STLS}} - \hat{x}_{\text{STLS}}\|_2 \leq \eta \frac{(1 + \|x_{\text{STLS}}\|_2)\sigma_{k+1}(\widehat{C}) + (\|x_{\text{STLS}}\|_2 + \lambda^{-1})\sigma_{k+1}(C)}{(\sigma_k(A) - \sigma_{k+1}(C) - 2\eta)\sigma_{k+1}(\widehat{C})},$$

provided that $\sigma_k(A) - \sigma_{k+1}(C) > 2\eta$.

This theorem shows that if the perturbation $\eta = \|\Delta C\|_2$ is small, we can expect a small error $\|x_{\text{STLS}} - \hat{x}_{\text{STLS}}\|_2$ as long as $\sigma_k(A)$ and $\sigma_{k+1}(C)$ are not very close.

5. Numerical experiments

In this section, we compare our bounds given by Theorems 3.1 and 4.1 with their existing counterparts given in [4] and [9]. The experiments were carried out in MATLAB. As we shall see, our numerical experiments have shown that our bounds are very close to the existing ones. It is probably very difficult, if possible, to prove which ones are superior. Moreover, the bounds in [9] are valid for the differences between the TLS and LS, whereas our bounds are valid for more general STLS and LS. Thus we compare the bounds in [9] with our bounds for the special case when $\lambda = 1$.

Our test matrix is the matrix of rank k that is closest in the matrix 2-norm to a randomly generated matrix with entries uniformly distributed on $[0, 1]$. Our test right-hand side vector is a vector with entries uniformly distributed on $[0, 1]$. The STLS solution is given by $x_{\text{STLS}} = -\lambda^{-1}V_{12}(v_{22}^T)^\dagger$ and the LS solution by $x_{\text{LS}} = A^\dagger b$.

Example 1. Theorem 3.1 gives two bounds for $\|x_{\text{STLS}} - x_{\text{LS}}\|_2$. As shown in the results, the second one is slightly larger than the first, but much simpler. In [9, (3.4)], Wei gives the bound, using our notations,

$$\|x_{\text{TLS}} - x_{\text{LS}}\|_2 \leq \rho^2 \sqrt{\|x_{\text{TLS}}\|_2^2 + 1} + \rho \|x_{\text{TLS}}\|_2.$$

To compare, we set $\lambda = 1$ in our bounds. Table 1 shows that our bounds are only slightly larger than Wei’s bound. However, Wei’s bound is implicit in that it involves the solution norm $\|x_{\text{TLS}}\|_2$, whereas our second bound is simpler and can be obtained without $\|x_{\text{STLS}}\|_2$. Moreover, our bounds are more general in that they are for STLS with TLS as a special case when $\lambda = 1$.

Example 2. Paige and Strakoš [4] give an expression for $\|\bar{r}_{\text{STLS}}\|_2$ [4, (4.13)] and an approximation for $\|r_{\text{LS}}\|_2$ [4, (4.14)], from which, using our notations, we have

Table 1

Comparison of the bound in [9] and the first and second bounds in Theorem 3.1

Size	Rank	$\ x_{\text{TLS}} - x_{\text{LS}}\ _2$	Bound in [9]	First bound	Second bound
200 × 100	80	1.4094	3.6853	3.7912	4.0397
400 × 300	100	0.2455	1.0248	1.2000	1.7626
800 × 700	600	1.3091	3.8888	4.1931	4.4009

Table 2

Comparison of the estimate for $\|\bar{r}_{\text{STLS}}\|_2 - \|r_{\text{LS}}\|_2$ from [4] and the bound in Theorem 3.1

λ	$\ \bar{r}_{\text{STLS}}\ _2 - \ r_{\text{LS}}\ _2$	Estimate from [4]	Bound in Theorem 3.1
0.005	1.111E−9	2.221E−9	2.302E−3
0.5	1.627E−1	2.816E−1	2.405E+0
1	1.605E+0	2.298E+0	4.703E+0
5	3.012E+0	4.169E+0	6.363E+0
10	3.196E+0	4.250E+0	6.454E+0

Table 3

Comparison of the bounds for perturbation in the TLS solution

Rank	$\ x_{\text{TLS}} - \hat{x}_{\text{TLS}}\ _2$	Bound in [9, (7.9)]	Bound in Theorem 4.1
500	4.787E−4	4.436E+1	4.254E−1
600	9.253E−5	7.083E+1	8.840E−2
800	2.569E−5	1.539E+2	2.770E−2

$$\|\bar{r}_{\text{STLS}}\|_2 - \|r_{\text{LS}}\|_2 \approx \lambda^{-1} \left| \sigma_{k+1}(C) \sqrt{1 + \|\lambda z\|_2^2} - \sigma_{k+1}(C) \sqrt{1 + \|\lambda x_{\text{LS}}\|_2^2} \right|,$$

where $z = (A^T A - \sigma_{k+1}^2(C)I)^{-1} (A^T b)$. To compare the above estimate with our bound in Theorem 3.1, we generated random A and b and used various values of λ . Table 2 shows the results for a 500×400 random matrix A with rank 350.

Table 2 shows that our bound for $\|\bar{r}_{\text{STLS}}\|_2 - \|r_{\text{LS}}\|_2$ given in Theorem 3.1 is about 1.5 times as large as the estimate from [4] when λ is larger than 1, which indicates that our bound is quite close for $\lambda > 1$. When λ is less than 1, our bound is larger, which means that our bound does not converge to zero as fast as the one from [4]. However, the evaluation of the expression for $\|\bar{r}_{\text{STLS}}\|_2$ in [4] involves solving for z in the system $(A^T A - \sigma_{k+1}^2(C)I)z = A^T b$. Whereas our bound for $\|\bar{r}_{\text{STLS}}\|_2 - \|r_{\text{LS}}\|_2$ given in Theorem 3.1 can be readily obtained.

Example 3. In this example, we compare our STLS solution perturbation bounds given in Theorem 4.1 for $\lambda = 1$ with the TLS solution perturbation bound [9, (7.9)]

$$\|x_{\text{TLS}} - \hat{x}_{\text{TLS}}\|_2 \leq \frac{6(\eta + \sigma_{k+1}(C))}{\sigma_k(A) - \sigma_{k+1}(C)} \sqrt{\|x_{\text{TLS}}\|_2^2 + 1}.$$

We generated random matrices A of size 1000×800 with various ranks k and random vectors b . Then we constructed random perturbation matrices ΔC such that $\|\Delta C\|_2 = \eta$, where η was set to $(\sigma_k(A) - \sigma_{k+1}([A, b]))/6$. A typical value of η was 0.0867.

The results in Table 3 show that our bound is tighter than [9, (7.9)]. This can be explained by the fact that our bound has η as a factor. For the same reason, our experiments also show that for small $\eta \approx \mathbf{c}\mathbf{u}\|C\|_2$, our bound gives close estimate for the perturbation in the solution.

6. Conclusion

In this paper, we first present quantitative relations between the scaled total least squares and least squares solutions and residuals. Theorem 3.1 shows that the two solutions and two residuals equal when $\lambda \rightarrow 0$. They can be very different when λ is not small and $\|v_{22}\|_2$ is small. Second, we give a perturbation analysis of the scaled total least squares problem. Theorem 4.1 shows that the solution of the perturbed problem is close to the original problem if the perturbation is small and $\sigma_{k+1}(C)$ and $\sigma_k(A)$ are not close to each other. Finally, our numerical experiments demonstrate that our bounds are competitive with existing bounds.

References

- [1] R.D. Fierro, J.R. Bunch, Bounding the subspaces from rank revealing two-sided orthogonal decompositions, *SIAM J. Matrix Anal. Appl.* 16 (1995) 743–759.
- [2] G.H. Golub, C.F. Van Loan, *Matrix Computations*, third ed., The Johns Hopkins University Press, Baltimore, MD, 1996.
- [3] S. Van Huffel, H. Zha, An efficient total least squares algorithm based on a rank-revealing two-sided orthogonal decomposition, *Numer. Algorithm* 4 (1993) 101–133.
- [4] C.C. Paige, Z. Strakoš, Bounds for the least squares distance using scaled total least squares, *Numer. Math.* 91 (2002) 93–115.
- [5] C.C. Paige, Z. Strakoš, Scaled total least squares fundamentals, *Numer. Math.* 91 (2002) 117–146.
- [6] B.D. Rao, Unified treatment of ls, tls and truncated svd methods using a weighted tls framework, in: S. Van Huffel (Ed.), *Recent Advances in Total Least Squares Techniques and Errors-in-Variables Modeling*, SIAM, Philadelphia, PA, 1997, pp. 11–20.
- [7] G.W. Stewart, J. Sun, *Matrix Perturbation Theory*, Academic Press, San Diego, CA, 1990.
- [8] G. Wang, Y. Wei, S. Qiao, *Generalized Inverses: Theory and Computations*, Science Press, Beijing, New York, 2004.
- [9] M. Wei, Algebraic relations between the total least squares and least squares problems with more than one solution, *Numer. Math.* 62 (1992) 123–148.
- [10] M. Wei, The analysis for the total least square problem with more than one solution, *SIAM J. Matrix Anal. Appl.* 13 (1992) 746–763.