

Article

Roundoff Error Analysis of an Algorithm Based on Householder Bidiagonalization for Total Least Squares Problems [†]

Zhanshan Yang ¹ and Xilan Liu ^{2,*}

¹ School of Mathematics and Statistics, Qinghai Nationalities University, Xining 810007, China; yangzhsh15@lzu.edu.cn

² School of Mathematics and Information Science, Baoji University of Arts and Sciences, Baoji 721000, China

* Correspondence: doclanliu@163.com

[†] Supported by the NNSF of China (11571004, 11701456), Natural Science Foundation of Qinghai Province (2018-ZJ-717), Foundation Sciences Qinghai Nationalities University (2020XJG11, 2019XJZ10).

Abstract: For large-scale problems, how to establish an algorithm with high accuracy and stability is particularly important. In this paper, the Householder bidiagonalization total least squares (HBITLS) algorithm and nonlinear iterative partial least squares for total least squares (NIPALS-TLS) algorithm were established, by which the same approximate TLS solutions was obtained. In addition, the propagation of the roundoff error for the process of the HBITLS algorithm was analyzed, and the mixed forward-backward stability of these two algorithms was proved. Furthermore, an upper bound of roundoff error was derived, which presents a more detailed and clearer approximation of the computed solution.



Citation: Yang, Z.; Liu, X. Roundoff Error Analysis of an Algorithm Based on Householder Bidiagonalization for Total Least Squares Problems.

Mathematics **2021**, *9*, 2550. <https://doi.org/10.3390/math9202550>

Academic Editors: Mehdi Salimi, Alicia Cordero Barbero and Christopher Goodrich

Received: 16 August 2021

Accepted: 30 September 2021

Published: 12 October 2021

Keywords: Householder bidiagonalization; NIPALS; roundoff error; total least squares problems

1. Introduction

Consider estimating x from the overdetermined linear system

$$Ax \approx b \quad \text{for } A \in \mathbb{R}^{m \times n}, \quad b \in \mathbb{R}^m \quad \text{and} \quad x \in \mathbb{R}^n, \quad (1)$$

where the error exists in both the right-hand side b and the data matrix A and $m \geq n + 1$. In this case, the total least squares (TLS) model should be appropriate to adopt (cf. [1,2]). The TLS approach is just to find a perturbation with the minimum Frobenius norm to make the system (1) a compatible system

$$\min \| (E, r) \|_F, \quad \text{subject to } b + r \in \text{Range}(A + E). \quad (2)$$

The TLS method is widely used in various scientific fields, such as physics, automatic control, signal processing, statistics, economics, biology, medicine etc. In essence, a solution of a TLS problem can be expressed by a singular value decomposition of the augmented matrix (A, b) . When the dimensions of A are not too large, one can use the truncated-SVD (TSVD) method. When the dimensions of A become large, this approach becomes prohibitive because the SVD algorithm is of complexity $\mathcal{O}(mn^2)$. The above considerations lead us to consider Krylov iterative methods, that do not alter the matrix A . The methods have the attractive feature just like the Lanczos methods—that when n increases, the computed extreme singular elements rapidly become good approximations to the exact ones, and are satisfactorily accurate even if k is far less than n theoretically [1]. Nevertheless, the orthonormal properties of the Krylov basis strongly support the use of these Householder matrix-based algorithms. This is particularly true when we need to be sure that the perturbed problem we are solving has to conserve some spectral similarity properties. This



Copyright: © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

will be especially relevant when we need to compute approximations of the TLS problems. In view of this, we consider applying the Householder bidiagonalization algorithm and the NIPALS PLS algorithm posed by Å. Björck [3] to TLS problems, the formed Householder bidiagonalization total least squares (HBITLS) algorithm, and NIPALS-TLS algorithm, respectively. Furthermore, we find that the HBITLS and NIPALS-TLS algorithms also compute the same approximate solutions for the TLS problems.

When it comes to practical problems, the arithmetic will be inaccurate and there will be errors in each step of the calculation. Arithmetic operations running on the computer have finite precision, so there will be rounding errors as long as there are numerical computations. These rounding errors cause the calculation quantities to be different from their theoretical values. One of the design principles of the floating-point operation is that it should encourage experts to develop robust, efficient, and portable numerical programs, enable the handling of arithmetic exceptions, and provide for the development of transcendental functions and high-precision arithmetic [4]. The results in the roundoff error analysis in Lanczos-type methods obtained by Paige [5–7] played an important role in interpreting the behavior of the Lanczos method in finite-precision computations. Parlett and Scott [8] used the results of the roundoff error analysis as the basis for suggesting a modification of the Lanczos method, which they called selective orthogonalization [8–10]. In addition, in many practical problems, the stop criterion can be safely selected on the basis of the rounding error analysis of the original problem, thereby diminishing the need for an extremely precise approximation of the algebraic problem solution [4]. As far as we know, the roundoff error analysis of the approximation TLS solutions obtained by using the Householder bidiagonalization procedure was not systematically performed in the literature. Hence, in this paper, we analyzed the propagation of the roundoff error during the process of the HBITLS algorithm and found that the HBITLS algorithm and NIPALS-TLS algorithm are mixed forward-backward stable.

The paper is organized as follows. The HBITLS algorithm and NIPALS-TLS algorithm were established, by which the same approximate TLS solution was obtained in Section 2. Section 3 analyzes the propagation of the roundoff error during the process of the HBITLS algorithm. A brief conclusion is shown in the last section.

2. HBITLS Algorithm and NIPALS-TLS Algorithm

It is well known that algorithms based on a sequence of orthogonal transformations with Householder matrices have very good stability properties; see Higham [4]. Based on this, this paper gives the HBITLS and NIPALS-TLS algorithms and finds that they both compute the same approximate TLS solutions.

2.1. HBITLS Algorithm

Let us first describe the Householder bidiagonalization process just as shown in [3]. However, in this paper, the process is used in the augmented matrix (\mathbf{b}, A) for the TLS problem. The idea is to compute orthogonal matrices $U \in \mathbb{R}^{m \times (n+1)}$ and $V \in \mathbb{R}^{(n+1) \times (n+1)}$, such that

$$U^T(\mathbf{b}, A)V = \left(\begin{array}{c|cccc} \beta_1 & \alpha_1 & & & \\ & \beta_2 & \alpha_2 & & \\ & & \beta_3 & \ddots & \\ & & & \ddots & \alpha_n \\ & & & & \beta_{n+1} \end{array} \right) = (\beta_1 \mathbf{e}_1, B_n) \equiv C_n. \quad (3)$$

$U = G_1 \cdots G_n = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{n+1}]$ and $V = H_1 \cdots H_n = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{n+1}]$ can be determined as a product of Householder matrices in each iteration. Generally, G_k introduces zero in the k th column, while H_k sets zero for the appropriate entries in the k th row. This can be done by an algorithm named Householder. Given the reason of space, and known to all, we omit it here, see algorithm 5.4.2 in [1] for details.

From the above process of Householder bidiagonalization, we know that V can be rewritten as $V = \begin{pmatrix} 1 & 0 \\ 0 & V_{n-1} \end{pmatrix}$, $V_{n-1} = [v_1, v_2, \dots, v_n]$.

In detail, from (3), we have

$$U^T \mathbf{b} = U^T(\mathbf{b}, A) V \mathbf{e}_1 = U^T(\mathbf{b}, A) \begin{pmatrix} 1 & 0 \\ 0 & V_{n-1} \end{pmatrix} \mathbf{e}_1 = C_n \mathbf{e}_1 = \beta_1 \mathbf{e}_1$$

and

$$A \mathbf{v}_i = \alpha_i \mathbf{u}_i + \beta_{i+1} \mathbf{u}_{i+1}, \text{ for } i = 1, \dots, n. \quad (4)$$

From (3), we have

$$(\mathbf{b}, A)^T U = \begin{pmatrix} 1 & 0 \\ 0 & V_{n-1} \end{pmatrix} C_n^T$$

and there comes

$$A^T \mathbf{u}_1 = \alpha_1 \mathbf{v}_1,$$

$$A^T \mathbf{u}_i = \alpha_i \mathbf{v}_i + \beta_{i+1} \mathbf{v}_{i+1}, \text{ for } i = 2, \dots, n-1. \quad (5)$$

Let $V_k = H_1 H_2 \dots H_k \begin{pmatrix} I_k \\ 0 \end{pmatrix}$, $U_k = G_k \dots G_2 G_1 \begin{pmatrix} I_k \\ 0 \end{pmatrix}$, and $C_k = (\beta_1 \mathbf{e}_1, B_k)$ is a leading principal submatrix of order $k+1$ of the final bidiagonal matrix C_n . As we all know, if the exact arithmetic is used, we have $U_{k+1}^T U_{k+1} = I$, $V_k^T V_k = I$. However, in any case, the previous equations remain within machine precision. Then (4) and (5) can be rewritten as

$$(\mathbf{b}, A) V_k = U_{k+1} C_k, \quad (6)$$

$$(\mathbf{b}, A)^T U_{k+1} = V_k C_k^T + \alpha_{k+1} \mathbf{v}_{k+1} \mathbf{e}_{k+1}^T = V_{k+1} \bar{C}_k^T, \quad (7)$$

where $\bar{C}_k = (C_k, \alpha_{k+1} \mathbf{e}_k) \in \mathbb{R}^{k \times (k+1)}$.

After performing the k -step Householder bidiagonalization iterations, the TLS problem can be reduced onto the subspace generated by U_{k+1} and V_k . Then the reduced TLS problem (also see [11]) is as follows

$$\min \|U_{k+1}^T ((\mathbf{b}, A) - (\hat{\mathbf{b}}_k, \hat{A}_k)) \begin{pmatrix} 1 & 0 \\ 0 & V_k \end{pmatrix}\|_F \quad \text{subject to } U_{k+1}^T \hat{A}_k V_k \mathbf{y} = U_{k+1}^T \hat{\mathbf{b}}_k, \quad (8)$$

or

$$\min \|(\beta_1 \mathbf{e}_1, B_k) - (\hat{\mathbf{e}}_k, \hat{B}_k)\|_F \quad \text{subject to } \hat{B}_k \mathbf{y} = \hat{\mathbf{e}}_k, \quad (9)$$

where $\mathbf{e}_1 = (1, 0, \dots, 0)^T$, and \hat{B}_k and $\hat{\mathbf{e}}_k$ are generally full. As in LSQR, seek an approximate TLS solution

$$\mathbf{x}_k = V_k \mathbf{y}_k \in \mathcal{K}_k(A^T A, A^T \mathbf{b}),$$

where $\mathcal{K}_k(B, \mathbf{y})$ denotes the Krylov subspace $\text{span}\{\mathbf{y}, B\mathbf{y}, \dots, B^{k-1}\mathbf{y}\}$.

Let the SVD of $(\beta_1 \mathbf{e}_1, B_k) = \bar{W}_k \bar{\Sigma}_k \bar{Z}_k^T$ and if let

$$\begin{pmatrix} \gamma_k \\ \mathbf{z}_k \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & V_k \end{pmatrix} (\bar{Z}_k \mathbf{e}_{k+1}) \quad (10)$$

with

$$\bar{Z}_k = \begin{pmatrix} \bar{Z}_1 & \bar{Z}_2 \end{pmatrix} = \begin{pmatrix} \bar{Z}_{11} & \bar{z}_{12} & 1 \\ \bar{Z}_{21} & \bar{z}_{22} & k \end{pmatrix}, \quad (11)$$

then the approximate TLS solution is given by

$$\mathbf{x}_k = -\frac{\bar{z}_k}{\gamma_k} \in \mathcal{K}_k(A^T A, A^T \mathbf{b}). \quad (12)$$

Note that we only need the last singular vector $\bar{Z}_k \mathbf{e}_{k+1}$ to compute \mathbf{x}_k . To this extent, summarizing the above process, we can get the Householder bidiagonalization TLS (HBITLS) algorithm as follows:

Remark 1. A variant of Algorithm 1 can also be given, in which the product of the Householder transformations applying to vectors are replaced by operations that can be performed concurrently, to a large extent. This variation gives an efficient method for developing parallelism in the case of parallel computing matrix vector products. In regard to this variation of Algorithm 1, one can refer to [12] and we omit it here.

Algorithm 1 HBITLS

- 1: Initialize: $C = (\mathbf{b}, A)$, $U := I$, $V := I$.
 - 2: for $j = 1, \dots, n+1$
 - $[\mathbf{v}, \beta] = \text{Householder}(C(j:m, j));$
 - $P = I_{m-j+1} - \beta \mathbf{v} \mathbf{v}^T;$
 - $C(j:m, j:n+1) = PC(j:m, j:n+1);$
 - $G = \text{diag}\{I_{j-1}, P\};$
 - $U = UG;$
 - $\mathbf{u}_j = U(:, j);$ if $j \leq n-1$
 - $[\mathbf{v}, \beta] = \text{Householder}(C(j, j+1:n+1));$
 - $Q = I_{n-j} - \beta \mathbf{v} \mathbf{v}^T;$
 - $C(j, j+1:n+1) = C(j, j+1:n+1)Q;$
 - $H = \text{diag}\{I_j, Q\};$
 - $V = VH;$
 - $\mathbf{v}_j = V(:, j).$
 - end
 - 3: Compute the last singular triplet for matrix C_k by employing the implicit zero-shift QR algorithm.
 - 4: the approximate TLS solution is given by $\mathbf{x}_k = -\frac{\bar{z}_k}{\gamma_k}$, see (10).
 - 5: end
-

Using the recursions (4) and (5), the following properties of $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_i\}$ and $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_i\}$ can be proved.

Lemma 1. The sets $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_i\}$ and $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_i\}$ generated by Algorithm 1 are the orthonormal basis of $\mathcal{K}_i(A^T A, A^T \mathbf{b})$ and $\mathcal{K}_i(AA^T, AA^T \mathbf{b})$ respectively.

Proof. As a result of the facts that $\beta_1 \mathbf{v}_1 = A^T \mathbf{b}$, $\alpha_1 \mathbf{u}_1 = A \mathbf{v}_1 = AA^T \mathbf{b} / \beta_1$, $H_1 \mathbf{v}_1 = \mathbf{e}_1$ and the process of Householder bidiagonalization, for $1 \leq k \leq i$, it's easy to know that $\mathbf{v}_j = H_1 \cdots H_k \mathbf{e}_j$ for $j = 1, \dots, k$, i.e., that

$$H_1 \cdots H_k = [\mathbf{v}_1, \dots, \mathbf{v}_k, \dots], \quad 1 \leq k \leq i. \quad (13)$$

Certainly $\mathcal{K}_1(A^T A, A^T \mathbf{b}) = \mathbb{R}(\mathbf{v}_1)$. It clearly holds if $i = 1$. Suppose for some $i > 1$ that the iteration has produced $V_i = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_i]$ with orthonormal columns such that

$$\mathcal{K}_i(A^T A, A^T \mathbf{b}) = \text{span}\{\mathbf{v}_1, (A^T A)\mathbf{v}_1, \dots, (A^T A)^{i-1}\mathbf{v}_1\} = \text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_i\}.$$

It is easy to see from (4) that

$$V_i^T A^T A V_i = B_i^T B_i$$

and we have $V_i^T \mathbf{r}_i = 0$, where $\mathbf{r}_i = A^T A \mathbf{v}_i - (\alpha_i^2 + \beta_{i+1}^2)\mathbf{v}_i - \alpha_{i-1}\beta_{i-1}\mathbf{v}_{i-1}$. If $\mathbf{r}_i \neq 0$, then $\mathbf{v}_{i+1} = \mathbf{r}_i / \|\mathbf{r}_i\|_2$ is orthogonal to $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_i$. It follows that $\mathbf{v}_{i+1} \notin \mathcal{K}_i(A^T A, A^T \mathbf{b})$ and

$$\mathbf{v}_{i+1} \in \text{span}\{A^T A \mathbf{v}_i, \mathbf{v}_i, \mathbf{v}_{i-1}\} \subseteq \mathcal{K}_{i+1}(A^T A, A^T \mathbf{b}).$$

Thus, $V_{i+1}^T V_{i+1}^T = I$ and

$$\text{span}(V_{i+1}) = \mathcal{K}_{i+1}(A^T A, A^T \mathbf{b}).$$

On the other hand, if $\mathbf{r}_i = 0$, then $A^T A V_i = V_i B_i^T B_i$. This says that $\text{span}(V_i) = \mathcal{K}_i(A^T A, A^T \mathbf{b})$ is invariant for $A^T A$ and the induction is complete. The proof of $\text{span}(U_i) = \mathcal{K}_i(AA^T, AA^T \mathbf{b})$ is in a similar way. \square

The Householder matrices H_i and G_i need not be formed explicitly. In other words, the matrices V_k and U_k can also remain in product form in the HBITLS algorithm. In floating-point operations, the Householder transformation does not have to worry too much about the loss of orthogonality.

2.2. The NIPALS-TLS Algorithm

For the NIPALS PLS algorithm, one can see in [3,13]. In this paper, we want to use it to solve the TLS problems and then form the NIPALS-TLS algorithm. We can find that the HBITLS algorithm and NIPALS-TLS algorithm generate the same sequences, orthonormal base V_k . From the uniqueness of this base, and combined with the relationship between the two algorithms, we conclude that the two algorithms generate the same numerical solution \mathbf{x}_k .

In [3], it tells us that we can set $A_0 = A$, $\mathbf{b}_0 = \mathbf{b}$, for $k = 1, 2, \dots$, we can produce sequences \mathbf{u}_k and \mathbf{v}_k according to the following form:

$$\mathbf{v}_k = A_{k-1}^T \mathbf{b}_{k-1} / \mu_k, \quad \mu_k = \|A_{k-1}^T \mathbf{b}_{k-1}\|_2, \quad (14)$$

$$\mathbf{p}_k = A_{k-1} \mathbf{v}_{k-1} / \rho_k, \quad \rho_k = \|A_{k-1} \mathbf{v}_{k-1}\|_2, \quad (15)$$

$$(A_k, \mathbf{b}_k) = (I - \mathbf{p}_k \mathbf{p}_k^T)(A_{k-1}, \mathbf{b}_{k-1}). \quad (16)$$

In (16) A_k and \mathbf{b}_k are formed by deflated A_{k-1} and \mathbf{b}_{k-1} by subtracting their orthogonal projections onto \mathbf{p}_k . We know that this operation uses elementary orthogonal transformations, such that $S = I - \mathbf{p} \mathbf{p}^T$, $\|\mathbf{p}\|_2 = 1$. The deflation in (16) can also be written as

$$A_k = A_{k-1} - \mathbf{p}_k \mathbf{s}_k^T, \quad \mathbf{s}_k = A_{k-1}^T \mathbf{p}_k, \quad (17)$$

$$\mathbf{b}_k = \mathbf{b}_{k-1} - \mathbf{p}_k \zeta_k^T, \quad \zeta_k = \mathbf{b}_{k-1}^T \mathbf{p}_k. \quad (18)$$

The process is terminated when it meets either $\|A_{k-1}^T \mathbf{b}_{k-1}\|_2 = 0$ or $\|A_{k-1} \mathbf{v}_i\|_2 = 0$. We note that if $\mathbf{p}_k^T A_{k-1} \mathbf{v}_k \neq 0$, then the rank of the matrix A_k is one less than that of A_{k-1} exactly.

Using exact arithmetic, the sets $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ and $\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_k\}$ generated by (14) and (15) are the unique orthogonal bases for the Krylov sub-spaces $\mathcal{K}_k(A^T A, A^T \mathbf{b})$ and $\mathcal{K}_k(AA^T, AA^T \mathbf{b})$, respectively. Summing (17) and (18) generates

$$A = P_k S_k^T + A_k, \quad \mathbf{b} = P_k \mathbf{z}_k + \mathbf{b}_k, \quad (19)$$

where $P_k = [p_1, p_2, \dots, p_k]$, $S_k = [s_1, s_2, \dots, s_k]$, and $z_k = (\zeta_1, \zeta_2, \dots, \zeta_k)^T$. These relationships maintain working accuracy and do not depend on orthogonality. The matrix $P_k S_k^T$ is a rank- k approximation to the data matrix A . From [3], we have $S_k^T = P_k^T A$ and $S_k^T V_k = R_k$. Thus, in exact arithmetic, the matrix $S_k^T V_k$ is upper bidiagonal with its elements

$$\theta_k = s_{k-1}^T v_k \quad \rho_k = s_k^T v_k = \|A_{k-1} v_k\|_2, \quad (20)$$

and

$$R_k \equiv \begin{pmatrix} \rho_1 & \theta_2 & & & \\ & \rho_2 & \theta_3 & & \\ & & \ddots & \ddots & \\ & & & \rho_{k-1} & \theta_k \\ & & & & \rho_k \end{pmatrix}.$$

By Paige [14], we know that R_k must be identical to the matrix that would be obtained from the conventional QR factorization of B_k , such that

$$Q_k B_k = \begin{pmatrix} R_k \\ 0 \end{pmatrix}.$$

Then we have

$$\begin{pmatrix} B_k & \beta_1 e_1 \end{pmatrix} = Q_k^T \begin{pmatrix} R_k & z_k \\ & \|b_k\|_2 \end{pmatrix} \equiv Q_k^T N_k. \quad (21)$$

Let $N_k = \check{U}_k \check{\Sigma}_k \check{V}_k$, and $\check{V}_k e_{k+1} \equiv \check{v}_{k+1} = \begin{pmatrix} (\check{v}_{k+1}^{(1)})^T & \check{v}_{k+1}^{(2)} \end{pmatrix}^T$, then the solution of the projected TLS problem (9) is $y_k = -\frac{\check{v}_{k+1}^{(1)}}{\check{v}_{k+1}^{(2)}}$, and the TLS solution is $x_k = -V_k \frac{\check{v}_{k+1}^{(1)}}{\check{v}_{k+1}^{(2)}}$. And the following theorem comes by (12) and (21)

Theorem 1. *the HBITLS and NIPALS-TLS algorithms compute the same approximate solutions x_k .*

3. Roundoff Error Analysis

In this section, we analyze the propagation of roundoff error during the process of the HBITLS algorithm and get the mixed forward-backward stability of the HBITLS algorithm and NIPALS-TLS algorithm naturally. The total roundoff error during the process of the HBITLS algorithm can be divided into the following four parts:

First, we can find that the HBITLS algorithm solves the original TLS problem (2) to a perturbed TLS problem. The propagation of the roundoff error of a Householder matrix in the HBITLS algorithm is advantageous when performing numerical computations.

From now on, we will denote by ε the machine precision under consideration. In [15], it shows that the computed Householder matrix $fl(H)$ comes near the exact Householder matrix H itself:

$$\|fl(H) - H\|_2 \leq 84\varepsilon + \mathcal{O}(\varepsilon^2).$$

Moreover, for a vector y , the computed updates with $fl(H)$ are very close to the exact updates with H :

$$fl(fl(H)y) = H(y + w), \quad \|w\|_2 \leq 87\varepsilon \|y\|_2 + \mathcal{O}(\varepsilon^2).$$

and, in general,

$$fl(fl(H_1) \cdots fl(H_j)y) = H_1 \cdots H_j(y + z), \quad \|z\|_2 \leq 87j\varepsilon \|y\|_2 + \mathcal{O}(\varepsilon^2). \quad (22)$$

The following lemma tells us that the reduced system calculated by the HBITLS algorithm is equivalent to the system formed after the original system has been disturbed. \bar{B}_k , \hat{Q}_k and \hat{P}_{k+1} are the floating-point computation of the matrices B_k , V_k and U_{k+1} in HBITLS algorithm, respectively.

Lemma 2. Let \bar{B}_k be the computed bidiagonalization matrix $(k+1) \times k$ matrix obtained by the HBITLS algorithm. Then, there comes a perturbation matrix E and exists two column orthogonal matrices \hat{Q}_k and \hat{P}_{k+1} s.t.

$$(A + E)\hat{Q}_k = \hat{P}_{k+1}\bar{B}_k,$$

and

$$\|E\|_2 \leq \sqrt{n}(174n + 3\sqrt{n} + 87)\varepsilon\|A\|_2 + \mathcal{O}(\varepsilon^2),$$

where n is the number of columns of matrix A . Furthermore, the matrix \hat{Q}_k is an orthonormal basis of $\mathcal{K}_k((A + E)^T(A + E), (A + E)^T(\mathbf{b} + \mathbf{e}))$ with a perturbation vector \mathbf{e} , where $\|\mathbf{e}\|_2 \leq 87\varepsilon\|\mathbf{b}\|_2 + \mathcal{O}(\varepsilon^2)$.

Proof. We prove this theorem by induction. The key point is that we should show the computed matrix, which will be shown by introduction from (3), for $G_k^T \cdots G_1^T[v_1, Av_1, \dots, Av_k]$ as follows

$$\hat{P}_{k+1} \cdots \hat{P}_1[\bar{v}_1 + \mathbf{g}_1, (A + G_1)\bar{v}_1 + \mathbf{w}_1, \dots, (A + G_k)\bar{v}_k + \mathbf{w}_k].$$

For $k = 1$, first, let $\mathbf{u}_1 = \mathbf{b}/\|\mathbf{b}\|_2$, $\bar{\mathbf{u}}_1 = fl(\mathbf{u}_1)$, a Householder matrix P_1 is found s.t. $P_1\mathbf{u}_1 = \mathbf{e}_1$. Set $\bar{P}_1 = fl(P_1)$, ref. [16] tells us that, corresponding to matrix \bar{P}_1 , we can find a Householder matrix to make

$$fl(\bar{P}_1\bar{v}_1) = \hat{P}_1(\bar{v}_1 + \mathbf{g}_1),$$

with

$$\|\mathbf{g}_1\|_2 \leq 87\varepsilon\|\bar{v}_1\|_2 + \mathcal{O}(\varepsilon^2) = 87\varepsilon + \mathcal{O}(\varepsilon^2).$$

Next, let $\mathbf{v}_1 = A^T\mathbf{u}_1/\|A^T\mathbf{u}_1\|_2$, $fl(A^T\mathbf{u}_1) = (A + G_0)^T\bar{\mathbf{u}}_1$, where $\|G_0\|_2 \leq \varepsilon 3\sqrt{n}\|A\|_2$. Similarly, for Av_1 , the computed result can be written as $fl(Av_1) = (A + G_1)\bar{v}_1$, where $\|G_1\|_2 \leq \varepsilon 3\sqrt{n}\|A\|_2$. Now, we set the Householder matrix P_2 s.t. $P_2P_1[v_1, Av_1]$ is upper bidiagonal matrix. We know P_2 only works the vector P_1Av_1 , so there's no change for the 1st column of $fl(\bar{P}_1[fl(v_1), fl(Av_1)])$ when produced by \bar{P}_2 and \hat{P}_2 . Likely, there's a Householder matrix \hat{P}_2 s.t. $\hat{P}_2fl(\bar{P}_1[\bar{v}_1, A\bar{v}_1])$ is bidiagonal matrix in theory, but the algorithm computes a matrix \bar{P}_2 in practice [16] such that

$$fl(\bar{P}_2\bar{P}_1fl(Av_1)) = \hat{P}_2\hat{P}_1((A + G_1)\bar{v}_1 + \mathbf{w}_1),$$

where

$$\|\mathbf{w}_1\|_2 \leq 174\varepsilon(\|A\|_2 + \|G_1\|_2)\|\bar{v}_1\|_2 = 174\varepsilon\|A\|_2 + \mathcal{O}(\varepsilon^2).$$

Finally, we have

$$\begin{aligned} fl(\bar{P}_2\bar{P}_1[\bar{v}_1, fl(Av_1)]) &= [fl(\bar{P}_1\bar{v}_1), fl(\bar{P}_2fl(\bar{P}_1fl(Av_1)))] \\ &= [\hat{P}_1(\bar{v}_1 + \mathbf{g}_1), \hat{P}_2\hat{P}_1((A + G_1)\bar{v}_1 + \mathbf{w}_1)] \\ &= \hat{P}_2\hat{P}_1[\bar{v}_1 + \mathbf{g}_1, (A + G_1)\bar{v}_1 + \mathbf{w}_1]. \end{aligned}$$

For the k th step, assume that the HBITLS algorithm has calculated the matrices $\bar{P}_1, \dots, \bar{P}_k$, associated with the Householder matrices $\hat{P}_1, \dots, \hat{P}_k$. Then, after k steps, we can get the following result:

$$\begin{aligned} fl(\bar{P}_k \cdots \bar{P}_1[\bar{v}_1, fl(Av_1), \dots, fl(Av_k)]) \\ = \hat{P}_k \cdots \hat{P}_1[\bar{v}_1 + \mathbf{g}_1, (A + G_1)\bar{v}_1 + \mathbf{w}_1, \dots, (A + G_k)\bar{v}_k + \mathbf{w}_k], \end{aligned}$$

where $\|w_i\|_2 \leq 87(i+1)\varepsilon\|A\|_2 + \mathcal{O}(\varepsilon^2)$.

We know $v_k = Q_1 \cdots Q_k e_k$, and $\bar{v}_k = fl(\bar{Q}_1 \cdots \bar{Q}_k e_k) = \hat{Q}_1 \cdots \hat{Q}_k (e_k + f_k)$, where $\|f_k\|_2 \leq 87k\varepsilon + \mathcal{O}(\varepsilon^2)$. For Av_k , the floating-point vector is $fl(Av_k) = (A + G_k)\bar{v}_k$, where $\|G_k\|_2 \leq \varepsilon 3\sqrt{n}\|A\|_2$. Likely, there is a Householder matrix P_{k+1} , which only works on the vector $P_k \cdots P_1 Av_k$, s.t. $P_{k+1}P_k \cdots P_1[v_1, Av_1, \dots, Av_k]$ is the upper bidiagonal matrix. The algorithm computes a matrix \bar{P}_{k+1} in practice so that

$$fl(\bar{P}_{k+1} \cdots \bar{P}_1 fl(Av_k)) = \hat{P}_{k+1} \cdots \hat{P}_1 ((A + G_k)\bar{v}_k + w_k),$$

where

$$\|w_k\|_2 \leq 87(k+1)\varepsilon(\|A\|_2 + \|G_k\|_2)(1 + \|f_k\|_2) \leq 87(k+1)\varepsilon\|A\|_2 + \mathcal{O}(\varepsilon^2).$$

Then the floating-point matrix is obtained, such that

$$\hat{P}_{k+1} \cdots \hat{P}_1 [\bar{v}_1 + g_1, (A + G_1)\bar{v}_1 + w_1, \dots, (A + G_k)\bar{v}_k + w_k].$$

Let $\hat{P}^{(k)}, \hat{Q}^{(k)}$ be the matrices, such that $(\hat{P}^{(k)})^T = \hat{P}_k \cdots \hat{P}_1$ and $(\hat{Q}^{(k)})^T = \hat{Q}_k \cdots \hat{Q}_1$, respectively, we find that the first $(i-1)$ rows of each \hat{Q}_i is e_1, \dots, e_{i-1} . Let $q_i^{(j)}$ be the i -th column of $\hat{Q}^{(j)}$, then the results are as follows

$$q_i^{(j)} = q_i^{(k)}, \quad \forall k \geq j, \quad i = 1, \dots, j.$$

Then there comes

$$(\hat{P}^{(k+1)})^T [\bar{v}_1 + g_1, (A + G_1)\bar{v}_1 + w_1, \dots, (A + G_k)\bar{v}_k + w_k],$$

it is an $n \times (k+1)$ upper bidiagonal matrix. And, $\forall j \leq k$, we obtain

$$\begin{aligned} (A + G_j)\bar{v}_j + w_j &= (A + G_j)\hat{Q}_1 \cdots \hat{Q}_j (e_j + f_j) + w_j \\ &= A\hat{Q}^{(j)} e_j + G_j \hat{Q}^{(j)} e_j + A\hat{Q}^{(j)} f_j + G_j \hat{Q}^{(j)} f_j + w_j \\ &= Aq_j^{(j)} + G_j q_j^{(j)} + A\hat{Q}^{(j)} f_j + G_j \hat{Q}^{(j)} f_j + w_j. \end{aligned}$$

Since $\forall k \leq j$, $q_j^{(j)}$ is the j -th column q_j of the matrix $\hat{Q}^{(j)}$, if we denote by $y_j = G_j q_j^{(j)} + A\hat{Q}^{(j)} f_j + G_j \hat{Q}^{(j)} f_j + w_j$, we have

$$(A + G_j)\bar{v}_j + w_j = Aq_j + y_j,$$

and so we can obtain

$$\begin{aligned} &(\hat{P}^{(k+1)})^T [\bar{v}_1 + g_1, (A + G_1)\bar{v}_1 + w_1, \dots, (A + G_k)\bar{v}_k + w_k] \\ &= (\hat{P}^{(k+1)})^T [\bar{v}_1 + g_1, Aq_1 + y_1, \dots, Aq_k + y_k] \\ &= (\hat{P}^{(k+1)})^T [(\bar{v}_1, Aq_1, \dots, Aq_k) + (g_1, y_1, \dots, y_k)]. \end{aligned}$$

If we cut off the first column of the matrix, we can set \tilde{B}_k with $n \times k$ such that $\tilde{B}_k = (\tilde{B}_k^T, 0)^T$, here \tilde{B}_k is an $(k+1) \times k$ upper bidiagonal matrix. If we denote $F_i = [y_1, \dots, y_i]$, $\forall i$, then $\tilde{B}_k = (\hat{P}^{(k+1)})^T (A[q_1, \dots, q_k] + F_k)$. In addition, $\forall j \leq k$, let \hat{Q}_j be the matrix made up of the first j columns of $\hat{Q}^{(k)}$, we obtain $\hat{P}^{(k+1)} \tilde{B}_k = A\hat{Q}^{(k)} + F_k$, and, from the structure of \tilde{B}_k , there comes $\hat{P}_{k+1} \tilde{B}_k = A\hat{Q}_k + F_k$. Then we can write $F_k = F_n \hat{Q}_n^T \hat{Q}_k$ owing to $\hat{Q}_n^T \hat{Q}_k = \begin{pmatrix} I_k \\ 0 \end{pmatrix}$, and so

$$\hat{P}_{k+1} \tilde{B}_k = (A + F_n \hat{Q}_n^T) \hat{Q}_k.$$

If $E = F_n \hat{Q}_n^T$, we can finish the proof of the first part of the lemma, because

$$\begin{aligned}\|E\|_2 &\leq \|E\|_F = \|F_n\|_F \\ &= \sqrt{\sum_{j=1}^n \|\mathbf{y}_j\|_2^2} \\ &\leq \sqrt{\sum_{j=1}^n \|G_j \mathbf{q}_j^{(j)} + A \hat{Q}^{(j)} \mathbf{f}_j + \mathbf{w}_j\|_2^2} + \mathcal{O}(\varepsilon^2) \\ &\leq \sqrt{n} \max_j \|G_j \mathbf{q}_j^{(j)} + A \hat{Q}^{(j)} \mathbf{f}_j + \mathbf{w}_j\|_2 + \mathcal{O}(\varepsilon^2) \\ &\leq \sqrt{n}(174n + 3\sqrt{n} + 87)\varepsilon \|A\|_2 + \mathcal{O}(\varepsilon^2).\end{aligned}$$

Finally, we prove that the subspace spanned by the columns of the matrix \hat{Q}_k is an orthogonal basis of a Krylov space. Let $\tilde{A} = A + E$, $\tilde{\mathbf{b}} = \|A^T \mathbf{b}\|_2 \hat{Q}_k \mathbf{e}_1$ and form $\mathcal{K}_k(\tilde{A}^T \tilde{A}, \tilde{A}^T \tilde{\mathbf{b}})$. We know that $\tilde{\mathbf{b}} = \|A^T \mathbf{b}\|_2 \hat{Q}_k \mathbf{e}_1 = \|A^T \mathbf{b}\|_2 \hat{Q}_1 \mathbf{e}_1$ and set $\mathbf{e} = \tilde{\mathbf{b}} - A^T \mathbf{b}$, then we have $\|\mathbf{e}\|_2 = \|A^T \mathbf{b}\|_2 \|(\hat{Q}_1 - Q_1) \mathbf{e}_1\|_2 \leq 87\varepsilon \|A^T \mathbf{b}\|_2 + \mathcal{O}(\varepsilon^2)$.

We still prove the rest of the theorem by induction; that is, to prove

$$(\tilde{A}^T \tilde{A})^i \tilde{A}^T \tilde{\mathbf{b}} = \hat{Q}_k \mathbf{r}_i, \quad \forall i \leq k-1,$$

where each vector \mathbf{r}_i has only the first $(i+1)$ components, which are different from zero.

For $i = 1$, we have

$$\tilde{A}^T \tilde{A} \tilde{A}^T \tilde{\mathbf{b}} = \tilde{A}^T \tilde{A} \|A^T \mathbf{b}\|_2 \hat{Q}_k \mathbf{e}_1 = \|A^T \mathbf{b}\|_2 \hat{Q}_k \tilde{B}_k^T \tilde{B}_k \mathbf{e}_1 = \hat{Q}_k \mathbf{r}_1,$$

since the last component of the vector $\tilde{B}_k^T \tilde{B}_k \mathbf{e}_1$ is zero, in addition, except for the first two components, the rest of the components of vector \mathbf{r}_1 are all zero.

Suppose for a given i the following relation is true,

$$(\tilde{A}^T \tilde{A})^{i+1} \tilde{A}^T \tilde{\mathbf{b}} = (\tilde{A}^T \tilde{A}) \hat{Q}_k \mathbf{r}_i = \hat{Q}_{k+1} \tilde{B}_k^T \tilde{B}_k \mathbf{r}_i,$$

and in the next step, we will show it is true for $i+1$.

From the inductive hypothesis, we know that only the first $(i+1)$ components of \mathbf{r}_i are not zero; therefore, the last component is zero of the vector $\tilde{B}_k^T \tilde{B}_k \mathbf{r}_i$ with $(i+2)$ non-zero elements. Then there comes a conclusion that

$$(\tilde{A}^T \tilde{A})^{i+1} \tilde{A}^T \tilde{\mathbf{b}} = \hat{Q}_k \mathbf{r}_{i+1}$$

and, hence, the lemma is proved. \square

Based on Lemma 2 and Algorithm 1, for the bidiagonalization matrix \tilde{B}_k obtained by the HBITLS algorithm, one can find an orthonormal matrix \hat{Q}_k s.t.

$$(A + E) \hat{Q}_k = \hat{P}_{k+1} \tilde{B}_k,$$

where \hat{Q}_k is just an orthogonal basis of $\mathcal{K}_k((A + E)^T(A + E), (A + E)^T(\mathbf{b} + \mathbf{e}))$. Based on this, we know that the first part of HBITLS, in exact arithmetic, gives the exact basis of the perturbed Krylov space $\mathcal{K}_k((A + E)^T(A + E), (A + E)^T(\mathbf{b} + \mathbf{e}))$. A perturbation bound for TLS solutions is given by Xie and Wei [17], see Lemma 3, which is related to the smallest singular value σ_{n+1} of (A, \mathbf{b}) , the TLS solution \mathbf{x} , and the residual $\mathbf{r} = \mathbf{b} - A\mathbf{x}$. Let $\hat{A} = A + \Delta A$ and $\hat{\mathbf{b}} = \mathbf{b} + \Delta \mathbf{b}$. Then, the unique solution of the perturbed TLS problem can be expressed as $\hat{\mathbf{x}}$. Denote $\kappa_{\mathbf{b}} = \frac{\|\mathbf{b}\|_2}{\|\mathbf{x}\|_2} \|B_{\lambda}^{-1} A^T\|_2$ and $\kappa_A = \frac{\|A\|_2}{\|\mathbf{x}\|_2} (\|\mathbf{r}\|_2 \|B_{\lambda}^{-1}\|_2 + \|x\|_2 \|B_{\lambda}^{-1} A^T\|_2)$ with $B_{\lambda} = A^T A - \sigma_{n+1}^2 I$. The perturbation bound is obtained under the genericity condition $\bar{\sigma}_n > \sigma_{n+1}$, where $\bar{\sigma}_n$ is the smallest singular value of A .

Lemma 3 ([17]). Consider the TLS problem (2) and assume that the genericity condition holds. If $\|(\Delta A, \Delta \mathbf{b})\|_F$ is sufficiently small, then we obtain that

$$\frac{\|\mathbf{x} - \hat{\mathbf{x}}\|_2}{\|\mathbf{x}\|_2} \lesssim \kappa_{\mathbf{b}} \frac{\|\Delta \mathbf{b}\|_2}{\|\mathbf{b}\|_2} + \kappa_A \frac{\|\Delta A\|_2}{\|A\|_2}.$$

Suppose that \mathbf{x} and $\hat{\mathbf{x}}$ are the exact TLS solutions of $A\mathbf{x} \approx \mathbf{b}$ and $(A + E)\hat{\mathbf{x}} \approx \mathbf{b} + \mathbf{e}$ respectively. The error introduced in this part of the HBITLS algorithm is the inherent error, so we can give $\|\mathbf{x} - \hat{\mathbf{x}}\|_2$ by Lemma 3 and Lemma 2 easily, see Theorem 2.

Secondly, let us consider the error between the TLS solution of the system $(A + E)\hat{\mathbf{x}} \approx \mathbf{b} + \mathbf{e}$ and the approximation solution $\hat{\mathbf{x}}_k = \hat{V}_k \hat{\mathbf{y}}_k$ of the system computed by the HBITLS algorithm at step k with the exact arithmetic, i.e., $\hat{\mathbf{y}}_k$ is the exact solution of the reduced TLS problem

$$\min \|(\beta_1 \mathbf{e}_1, \bar{B}_k) - (\check{\mathbf{e}}_k, \check{B}_k)\|_F \quad \text{subject to} \quad \check{B}_k \hat{\mathbf{y}}_k = \check{\mathbf{e}}_k, \quad (23)$$

For convenience, define $\bar{V}_k = \begin{pmatrix} 1 & 0 \\ 0 & V_k \end{pmatrix}$ and let

$$(\mathbf{b}, A) = W \Sigma Z^T, \quad (24)$$

where W and $Z = [z_1, z_2, \dots, z_{n+1}]$ are orthogonal matrices of dimension m and $n + 1$, respectively, Σ is an $m \times (n + 1)$ diagonal matrix whose diagonal entries σ_i are the singular values of (\mathbf{b}, A) , sorted in non-increasing order. Let $\theta(\mathbf{u}, \mathbf{v})$ be the subspace angle between $\mathbb{R}(\mathbf{u})$ and $\mathbb{R}(\mathbf{v})$.

Lemma 4. Let \mathbf{x} denotes the essential TLS solution to the linear system (2) satisfying genericity condition and \mathbf{x}_k be the approximation solution obtained from Algorithm 1. Then

$$\sin \theta(\mathbf{z}, \mathbf{z}^{(k)}) \leq \|\mathbf{x} - \mathbf{x}_k\| \leq \sin \theta(\mathbf{z}, \mathbf{z}^{(k)}) \sqrt{1 + \|\mathbf{x}\|^2} \sqrt{1 + \|\mathbf{x}_k\|^2}, \quad (25)$$

where

$$\begin{aligned} \mathbf{z} &\equiv \mathbf{z}_{n+1} = \begin{pmatrix} z_{12} \\ z_{22} \end{pmatrix} \in \mathbb{R}^{n+1}, \\ \mathbf{z}^{(k)} &\equiv \begin{pmatrix} z_{12}^{(k)} \\ z_{22}^{(k)} \end{pmatrix} = \bar{V}_k \bar{\mathbf{z}}_{k+1} = \begin{pmatrix} 1 & 0 \\ 0 & V_k \end{pmatrix} \begin{pmatrix} \bar{z}_{12} \\ \bar{z}_{22} \end{pmatrix} \in \mathbb{R}^{n+1}. \end{aligned} \quad (26)$$

Proof. It is easy to know that

$$\begin{pmatrix} -1 \\ \mathbf{x} \end{pmatrix} - \begin{pmatrix} -1 \\ \mathbf{x}_k \end{pmatrix} = \begin{pmatrix} z_{12}^{(k)} \\ z_{22}^{(k)} \end{pmatrix} (z_{12}^{(k)})^{-1} - \begin{pmatrix} z_{12} \\ z_{22} \end{pmatrix} z_{12}^{-1}. \quad (27)$$

Then we have an orthonormal matrix $G \in \mathbb{R}^{(n+1) \times n}$ with the partition

$$G = \begin{bmatrix} G_1 \\ G_2 \end{bmatrix} \begin{matrix} n \\ 1 \end{matrix}$$

such that $G^T ((z_{12}^{(k)})^T (z_{22}^{(k)})^T)^T = 0$ (i.e., G “forms a complete space”). From Equation (27) there comes

$$G^T \left[\begin{pmatrix} -1 \\ \mathbf{x} \end{pmatrix} - \begin{pmatrix} -1 \\ \mathbf{x}_k \end{pmatrix} \right] = -G^T \begin{pmatrix} z_{12} \\ z_{22} \end{pmatrix} z_{12}^{-1}, \quad (28)$$

and, therefore

$$G_2^T(\mathbf{x} - \mathbf{x}_k) = -G^T \begin{pmatrix} z_{12} \\ z_{22} \end{pmatrix} z_{12}^{-1}. \quad (29)$$

From the CS theorem [18], we know that $\sigma_{\min}(G_2)^{-1} = (z_{12}^{(k)})^{-1}$. Then

$$\begin{aligned} \sigma_{\min}(G_2) \|\mathbf{x} - \mathbf{x}_k\| &\leq \|G_2^T(\mathbf{x} - \mathbf{x}_k)\| \\ &= \|G^T \begin{pmatrix} z_{12} \\ z_{22} \end{pmatrix} z_{12}^{-1}\| \\ &\leq \sin \theta(\mathbf{z}, \mathbf{z}^{(k)}) z_{12}^{-1}. \end{aligned}$$

It was noticed that

$$\sin \theta(\mathbf{z}, \mathbf{z}^{(k)}) = \|G^T \begin{pmatrix} z_{12} \\ z_{22} \end{pmatrix}\|$$

denotes the sine of the subspace angle between $\mathbb{R}(\mathbf{z})$ and $\mathbb{R}(\mathbf{z}^{(k)})$. Hence, the upper bound can be proved as follows

$$\begin{aligned} \|\mathbf{x} - \mathbf{x}_k\| &\leq \sin \theta(\mathbf{z}, \mathbf{z}^{(k)}) z_{12}^{-1} (z_{12}^{(k)})^{-1} \\ &= \sin \theta(\mathbf{z}, \mathbf{z}^{(k)}) \sqrt{1 + \|\mathbf{x}\|^2} \sqrt{1 + \|\mathbf{x}_k\|^2}. \end{aligned}$$

For the lower bound, we have

$$\begin{aligned} \|\mathbf{x} - \mathbf{x}_k\| &\geq \|G_2^T(\mathbf{x} - \mathbf{x}_k)\| \\ &= \|G^T \begin{pmatrix} z_{12} \\ z_{22} \end{pmatrix} z_{12}^{-1}\| \\ &\geq z_{12}^{-1} \sin \theta(\mathbf{z}, \mathbf{z}^{(k)}) \\ &\geq \sin \theta(\mathbf{z}, \mathbf{z}^{(k)}). \end{aligned}$$

Since $z_{12}^{-1} \geq 1$, and this proves the upper bound case. \square

Thirdly, we need to consider how to solve problem (9) and show that the error is between the solution obtained by this method and the theoretical solution. Let the computed solution be

- $\bar{\mathbf{x}}_k = \hat{V}_k \bar{\mathbf{y}}_k$, where $\bar{\mathbf{y}}_k$ is the computed solution of the problem (23).

In [19], James and Kahan posed an algorithm named QR iteration with a zero shift, which guaranteed forward stability. Furthermore, an implicit algorithm about it is given. Error analysis including the singular values and singular vectors are also given, which is just what we've needed.

Lemma 5 ([19]). *Let the matrix \bar{B} obtained by running the implicit zero-shift QR algorithm on a bidiagonal matrix B with $n \times n$. Suppose that all perturbation angles θ emerged from the operations of the algorithm satisfy $\sin^2 \theta \leq \tau < 1$. Let $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$ and $\bar{\sigma}_1 \geq \bar{\sigma}_2 \geq \dots \geq \bar{\sigma}_n$ are the singular values of B and \bar{B} respectively. If*

$$\omega \equiv \frac{88n\epsilon}{(1-\tau)^2} < 1, \quad (30)$$

then we have:

$$|\sigma_i - \bar{\sigma}_i| \leq \frac{\omega}{1-\omega} \sigma_i.$$

Moreover, let $\sigma_{k1} \geq \sigma_{k2} \geq \dots \geq \sigma_{kn}$ be the singular values of B_k produced after k steps of the implicit zero-shift QR algorithm. Then if condition (30) holds, and all perturbation angles θ satisfy $\sin^2 \theta \leq \tau < 1$, we obtain

$$|\sigma_i - \sigma_{ki}| \leq \left(\frac{1}{(1 - \omega)^k - 1} \right) \sigma_i \approx \frac{88kn\epsilon}{(1 - \tau)^2} \sigma_i.$$

James and Kahan [19] also give the relative differences between the singular vectors of B and the ones of \tilde{B} .

Lemma 6 ([19]). Let σ_i be the singular value of be an unreduced bidiagonal matrix B with \mathbf{u}_i and \mathbf{v}_i being its corresponding left and right singular vectors, respectively. Let $\tilde{\mathbf{u}}_i$ and $\tilde{\mathbf{v}}_i$ be the singular vectors computed by the implicit zero-shift QR algorithm. Then the bound of the errors in $\tilde{\mathbf{v}}_i$ are shown by

$$\theta(\tilde{\mathbf{v}}_i, \mathbf{v}_i) \leq p(n)\epsilon / rel_{gap} \equiv p(n)\epsilon / \min(|\sigma_i - \sigma_{i+1}| / \sigma_i). \quad (31)$$

Then, combining with the perturbation bound of TLS given in [20] as shown in Lemma 7, we can give the error estimate $\|\hat{\mathbf{y}}_k - \tilde{\mathbf{y}}_k\|_2$.

If let $(\tilde{A}, \tilde{\mathbf{b}})$ is a rank- k matrix approximation to (A, \mathbf{b}) , and $(\Delta\tilde{A}, \Delta\tilde{\mathbf{b}}) = (A, \mathbf{b}) - (\tilde{A}, \tilde{\mathbf{b}})$. Let $(\bar{A}, \bar{\mathbf{b}}) = (A, \mathbf{b}) + (\Delta A, \Delta \mathbf{b})$ represent a perturbation of (A, \mathbf{b}) , $(\check{A}, \check{\mathbf{b}})$ denote a rank- k matrix approximation to $(\bar{A}, \bar{\mathbf{b}})$ and define $(\Delta\check{A}, \Delta\check{\mathbf{b}}) = (\bar{A}, \bar{\mathbf{b}}) - (\check{A}, \check{\mathbf{b}})$, then

Lemma 7 ([20]). Let \mathbf{x} and $\tilde{\mathbf{x}}$ denote the TLS solution and the perturbed TLS solution. If $\max(\|\Delta\check{A}\|, \|\Delta\check{\mathbf{b}}\| + \|\Delta A\|) < \sigma'_k$ (the k -th singular value of A) may be provided. Then

$$\sin \theta(\mathbf{v}, \tilde{\mathbf{v}}) \leq \|\mathbf{x} - \tilde{\mathbf{x}}\| \leq \sin \theta(\mathbf{v}, \tilde{\mathbf{v}}) \sqrt{1 + \|\mathbf{x}\|^2} \sqrt{1 + \|\tilde{\mathbf{x}}\|^2}, \quad (32)$$

where \mathbf{v} and $\tilde{\mathbf{v}}$ are the smallest right singular vectors of (A, \mathbf{b}) and $(\bar{A}, \bar{\mathbf{b}})$ respectively.

In summary, if let $\tilde{\mathbf{x}}_k = fl(\tilde{V}_k \tilde{\mathbf{y}}_k)$ be the final computed solution at the k -th step, then roundoff error analysis of HBITLS algorithm for TLS problem can be shown as follows

Theorem 2. Considering the HBITLS algorithm at step k , the roundoff error emerged during the algorithm can be bounded as follows:

$$\begin{aligned} \|\mathbf{x} - \tilde{\mathbf{x}}_k\| &\leq \kappa_b 87\epsilon + \kappa_A \sqrt{n}(174n + 3\sqrt{n} + 87)\epsilon \\ &+ \sin \theta(\hat{\mathbf{z}}, \hat{\mathbf{z}}^{(k)}) \sqrt{1 + \|\hat{\mathbf{x}}\|^2} \sqrt{1 + \|\hat{\mathbf{x}}_k\|^2} \\ &+ \sin \theta(\hat{\mathbf{z}}^{(k)}, \hat{\mathbf{z}}) \sqrt{1 + \|\hat{\mathbf{y}}_k\|^2} \sqrt{1 + \|\hat{\mathbf{y}}_k\|^2} \\ &+ 87k\epsilon \|\hat{\mathbf{y}}_k\| + \mathcal{O}(\epsilon^2), \end{aligned}$$

where $\hat{\mathbf{z}}$ and $\hat{\mathbf{z}}^{(k)}$ are defined in (26) similarly, $\hat{\mathbf{z}}$ is the computed smallest right singular vector of $(\beta_1 \mathbf{e}_1, \tilde{B}_k)$.

Proof. The roundoff error can be composed of the following parts

$$\mathbf{x} - \tilde{\mathbf{x}}_k = (\mathbf{x} - \hat{\mathbf{x}}) + (\hat{\mathbf{x}} - \hat{\mathbf{x}}_k) + (\hat{\mathbf{x}}_k - \hat{V}_k \hat{\mathbf{y}}_k) + (\hat{V}_k \hat{\mathbf{y}}_k - fl(\tilde{V}_k \tilde{\mathbf{y}}_k))$$

and we analyze these errors separately.

For the first part, \mathbf{x} and $\hat{\mathbf{x}}$ are the TLS solutions of the systems $A\mathbf{x} \approx \mathbf{b}$ and $(A + E)\hat{\mathbf{x}} \approx \mathbf{b} + \mathbf{e}$, respectively, in line with Lemma 2, so the error of this part is the inherent error. Then, combining with Lemma 3, we have

$$\|\mathbf{x} - \hat{\mathbf{x}}\| \leq \kappa_b 87\epsilon + \kappa_A \sqrt{n}(174n + 3\sqrt{n} + 87)\epsilon,$$

where κ_A and κ_b , see Lemma 3.

For the second part, this error is owing to the approximate solution of $(A + E)\hat{x} \approx b + e$ obtained by using HBITLS algorithm after k steps with the exact arithmetic. Lemma 4 tells us that

$$\|\hat{x} - \hat{x}_k\| \leq \sin \theta(\hat{z}, \hat{z}^{(k)}) \sqrt{1 + \|\hat{x}\|^2} \sqrt{1 + \|\hat{x}_k\|^2}.$$

For the third part, it is noticed that $\hat{x}_k = \hat{V}_k \hat{y}_k$, we have that $\|\hat{x}_k - \hat{V}_k \bar{y}_k\| = \|\hat{y}_k - \bar{y}_k\|$, where $\|\hat{y}_k - \bar{y}_k\|$ is the roundoff error stem from the projected TLS solution. Since B_k is a special form of bidiagonal matrices, we consider using the implicit zero-shift QR algorithm to perform singular value decomposition. (31) gives an upper bound of the angle between the solution vectors, and combining Lemma 7, we know

$$\|\hat{x}_k - \hat{V}_k \bar{y}_k\| = \|\hat{y}_k - \bar{y}_k\| \leq \sin \theta(\hat{z}^{(k)}, \bar{z}) \sqrt{1 + \|\hat{y}_k\|^2} \sqrt{1 + \|\bar{y}_k\|^2},$$

where $\theta(\hat{z}^{(k)}, \bar{z})$ is the subspace angle between the sub-spaces produced by the smallest right singular vector and the computed smallest right singular vector of $(\beta_1 e_1, \bar{B}_k)$, respectively.

For the last part, we know $V_k = H_1 H_2 \cdots H_k \begin{pmatrix} I_k \\ 0 \end{pmatrix}$, where $H_1 H_2 \cdots H_k$ is the product of k Householder matrices. So, on the basis of (22), we obtain

$$\|\hat{V}_k \bar{y}_k - f_l(\bar{V}_k \bar{y}_k)\| \leq 87k\epsilon \|\bar{y}_k\| + \mathcal{O}(\epsilon^2).$$

□

By theorem 2, we get the mixed forward–backward stability of the HBITLS algorithm and NIPALS-TLS algorithm naturally. The backward stability will generate perturbation that will marginally influence the theoretical convergence of the residual to zero.

Remark 2. The bound we introduced in Theorem 2 shows that the total roundoff errors are dominated by the approximation errors $\|\hat{x} - \hat{x}_k\|$. From this, we can know that, in many practical problems, we can safely select the stopping criteria required by the algorithm based on the theoretical nature of the original problem. This shows that, in a great deal of practical studies, the stopping criteria may be effectively selected based on the theoretical properties of the problem itself, thereby reducing the cost required to pursue an extremely accurate approximate solution to the original problem.

4. Conclusions

For large-scale problems, how to give an algorithm with good accuracy and stability is particularly important. In this paper, the Householder bidiagonalization total least squares (HBITLS) algorithm and nonlinear iterative partial least squares (NIPALS-TLS) algorithm are given. The HBITLS uses the Householder bidiagonalization algorithm for reducing (b, A) to upper bidiagonal form and then runs the implicit zero-shift QR algorithm to compute the smallest right singular vector of the reduced form for the approximation solutions. The NIPALS-TLS is based on rank-reducing orthogonal projections. The two algorithms compute the same approximate TLS solutions. By analyzing the propagation of the roundoff error during the process of the HBITLS algorithm, we find that the HBITLS algorithm and the NIPALS-TLS algorithm are to be mixed forward–backward stable. In addition, in many practical problems, the stop criterion can be safely selected on the basis of the rounding error analysis of the original problem. The upper bound of our roundoff error gives a more detailed and clearer approximation of the computed solution.

Author Contributions: Conceptualization, Z.Y. and X.L.; methodology, Z.Y. and X.L.; formal analysis, Z.Y.; data curation, X.L.; writing—original draft preparation, X.L.; writing—review and editing, Z.Y. and X.L. All authors have read and agreed to the published version of the manuscript.

Funding: This research is supported by the NNSF of China (11571004, 11701456), Natural Science Foundation of Qinghai Province (2018-ZJ-717), Foundation Sciences Qinghai Nationalities University(2020XJG11, 2019XJZ10), Innovation team of Qinghai Nationalities University.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Golub, G.H.; Loan, C.F.V. *Matrix Computations*, 4th ed.; The John Hopkins University Press: Baltimore, MD, USA, 2013.
2. Huffel, S.V.; Vandwalle, J. The total least squares problem, computational aspects and analysis. In *Frontiers in Applied Mathematics*; SIAM: Philadelphia, PA, USA, 1991; Volume 9.
3. Björck, Å. Stability of two direct methods for bidiagonalization and partial least squares. *SIAM J. Matrix Anal. Appl.* **2014**, *35*, 279–291.
4. Higham, N.J. *Accuracy and Stability of Numerical Algorithms*, 2nd ed.; SIAM Press: Philadelphia, PA, USA, 2002.
5. Paige, C.C. The Computation of Eigenvalues and Eigenvectors of Very Large Sparse Matrices. Ph.D. Thesis, University of London, London, UK, 1971.
6. Paige, C.C. Computational variants of the Lanczos method for the eigenproblem. *J. Inst. Math. Appl.* **1972**, *10*, 373–381. [[CrossRef](#)]
7. Paige, C.C. Error analysis of the Lanczos algorithm for tridiagonalizing a symmetric matrix. *J. Inst. Math. Appl.* **1976**, *18*, 341–349. [[CrossRef](#)]
8. Parlett, B.N.; Scott, D.S. The Lanczos algorithm with selective orthogonalization. *Math. Comput.* **1979**, *33*, 217–238. [[CrossRef](#)]
9. Ikramov, K.D. Sparse matrices. *Itogi Nauki Tekh. Mat. Anal.* **1982**, *20*, 189–259. (In Russian)
10. Parlett, B.N. *The Symmetric Eigenvalue Problem*; Prentice-Hall: Englewood Cliffs, NJ, USA, 1980.
11. Björck, Å. *Numerical Methods for Least Squares Problems*; SIAM: Philadelphia, PA, USA, 1996.
12. Walker, H.F. Implementation of the GMRES method using Householder transformations. *SIAM J. Sci. Statist. Comput.* **1988**, *9*, 152–163. [[CrossRef](#)]
13. Wold, H. Estimation of principal components and related models by iterative least squares. In *Multivariate Analysis*; Krishnaiah, P.R., Ed.; Academic Press: New York, NY, USA, 1966; pp. 391–420.
14. Paige, C.C.; Saunders, M.A. LSQR: An algorithm for sparse linear equations and sparse least squares. *ACM Trans. Math. Software* **1982**, *8*, 43–71. [[CrossRef](#)]
15. Wilkinson, J.H. *The Algebraic Eigenvalue Problem*; Oxford University Press: London, UK, 1965.
16. Lawson, C.; Hanson, R. *Solving Least Squares Problems*; Prentice Hall: Englewood Cliffs, NJ, USA, 1974.
17. Xie, P.; Xiang, H.; Wei, Y. A contribution to perturbation analysis for total least squares problems. *Numer. Algorithm.* **2017**, *75*, 381–395. [[CrossRef](#)]
18. Paige, C.C.; Saunders, M.A. Towards a generalized singular value decomposition. *SIAM J. Numer. Anal.* **1981**, *18*, 398–405. [[CrossRef](#)]
19. Demmel, J.; Kahan, W. Accurate singular values of bidiagonal matrices. *SIAM J. Sci. Statist. Comput.* **1990**, *11*, 873–912. [[CrossRef](#)]
20. Fierro, R.D.; Bunch, J.R. Perturbation theory for orthogonal projection methods with application to least squares and total least squares. *Linear Algebra Appl.* **1996**, *234*, 71–96. [[CrossRef](#)]