

# The exact condition number of the truncated singular value solution of a linear ill-posed problem

El Houcine Bergou, Serge Gratton, Jean Tshimanga Ilunga

## ▶ To cite this version:

El Houcine Bergou, Serge Gratton, Jean Tshimanga Ilunga. The exact condition number of the truncated singular value solution of a linear ill-posed problem. SIAM Journal on Mathematical Analysis, 2014, 35 (3), pp.1073-1085. 10.1137/120869286 . hal-02147972

# HAL Id: hal-02147972 https://hal.science/hal-02147972

Submitted on 5 Jun 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



## **Open Archive Toulouse Archive Ouverte**

OATAO is an open access repository that collects the work of Toulouse researchers and makes it freely available over the web where possible

This is a publisher's version published in: http://oatao.univ-toulouse.fr/22597

### **Official URL**

DOI : https://doi.org/10.1137/120869286

**To cite this version:** Bergou, El Houcine and Gratton, Serge and Tshimanga Ilunga, Jean *The exact condition number of the truncated singular value solution of a linear ill-posed problem.* (2014) SIAM Journal on Mathematical Analysis, 35 (3). 1073-1085. ISSN 0036-1410

### THE EXACT CONDITION NUMBER OF THE TRUNCATED SINGULAR VALUE SOLUTION OF A LINEAR ILL-POSED PROBLEM\*

#### EL HOUCINE BERGOU<sup> $\dagger$ </sup>, SERGE GRATTON<sup> $\dagger$ </sup>, AND JEAN TSHIMANGA<sup> $\dagger$ </sup>

Abstract. The main result of this paper is the formulation of an explicit expression for the condition number of the truncated least squares solution of Ax = b. This expression is given in terms of the singular values of A and the Fourier coefficients of b. The result is derived using the notion of the Fréchet derivative together with the product norm on the data [A, b] and the 2-norm on the solution. Numerical experiments are given to confirm our results by comparing them to those obtained by means of a finite difference approach.

Key words. truncated singular value decomposition, condition number, Fréchet derivative, least squares solution, perturbation theory

### DOI. 10.1137/120869286

1. Introduction. Perturbation analysis is the study of the sensitivity of the solution of a given problem to perturbations in the data. The concept of condition number allows us to assess the sensitivity of the solution. Sensitivity and conditioning theory has been applied to many fundamental problems of linear algebra, such as linear systems, linear least squares, or eigenvalue problems [1, 2, 8, 4, 13]. In this paper, we investigate the condition number for the so-called truncated singular value decomposition (TSVD) solution to linear least squares problems. TSVD solutions arise in a wide variety of applications in science, technology, and engineering. In inverse problems, for example, the TSVD can be considered as a regularization technique for ill-conditioned matrices with well-determined numerical rank; see [5, 6, 14]. Applications of TSVD solutions in this area include computational tomography, image deblurring, digital signal processing, and geophysical inversion in seismology. Some numerical solutions of partial differential equations may also require techniques such as TSVD; see [11].

Let A be an  $n \times p$  matrix  $(n \ge p)$  with rank $(A) = r^* \le p$  and let

$$A = U\Sigma V^{T}$$

be the *full* singular value decomposition of A with singular values of A arranged in descending order in  $\Sigma$ . Then, given an *n*-vector b, the least squares problem

$$\min_{x \in \Re p} \|Ax - b\|_2$$

has the minimum 2-norm solution  $x^* = V_{r^*} \Sigma_{r^*}^{-1} U_{r^*}^T b$ , where  $\Sigma_{r^*}$  is the diagonal matrix consisting of the first  $r^*$  singular values of A in descending order, and  $U_{r^*}$  and  $V_{r^*}$  are

<sup>\*</sup>This research was supported by the "Assimilation de Données pour la Terre, l'Atmosphère et l'Oc éa n (ADTAO)" project, funded by the Fondation "Sciences et Technologies pour l'Aéoro-nautique et l 'Espace (STAE)", Toulouse, France, within the "Réseau Thématique de Recherche Avancée (RTRA) ".

http://www.siam.org/journals/simax/35-3/86928.html

<sup>&</sup>lt;sup>†</sup>CERFACS, 42, avenue Gaspard Coriolis, 31057 Toulouse, France (elhoucine.bergou@cerfacs.fr, serge.gratton@cerfacs.fr, jean.tshimanga@cerfacs.fr).

formed from the first  $r^*$  columns of U and V, respectively. In some applications (e.g., problems arising from the discretization of an ill-posed problem), a better solution, in the sense that it is less sensitive than the original one to errors in the data (A, b), is obtained by a truncated least squares solution of the form

$$x_r = V_r \Sigma_r^{-1} U_r^T b,$$

for some  $r < r^*$ , and where  $V_r$ ,  $\Sigma_r$ , and  $U_r$  are defined as before but with r replacing  $r^*$ . It turns out that if  $\hat{U}_r$  and  $\hat{V}_r$  are any orthonormal bases for range  $(U_r)$  and range  $(V_r)$ , then

$$x_r = \hat{V}_r (\hat{U}_r^T A \hat{V}_r)^{-1} \hat{U}_r^T b.$$

Now let A and b be perturbed to yield  $\tilde{A} = A + E$  and  $\tilde{b} = b + f$ , and let  $\tilde{U}_r$  and  $\tilde{V}_r$  form a pair of bases for the left and right singular subspaces associated with the r first singular values of  $\tilde{A}$ . The corresponding truncated least squares solution of the perturbed problem is then

(1.1) 
$$\tilde{x}_r = \tilde{V}_r (\tilde{U}_r^T \tilde{A} \tilde{V}_r)^{-1} \tilde{U}_r^T b.$$

Now it turns out that if the Fréchet derivative,  $x'_r$ , of the function  $x_r$  exists, then we have

$$\tilde{x}_r = x_r + x'_r \cdot (E, f) + o(||(E, f)||).$$

Here,  $x'_r(E, f)$  is the application of a linear operator to (E, f). Given a norm on (E, f), call it  $\|.\|_{(\alpha,\beta)}$ , the *condition number* of  $x_r$  is defined to be the operator norm

$$|||x_r'|||_{(\alpha,\beta),2} = \max_{[\alpha E,\beta f] \neq 0} \frac{||x_r' \cdot (E,f)||_2}{||[E,f]||_{(\alpha,\beta)}}$$

The particular norm we use is defined by

$$||(E,f)||_{(\alpha,\beta)} = \sqrt{\alpha^2 ||E||_F^2 + \beta^2 ||f||_F^2},$$

where  $\|.\|_F$  is the usual Frobenius norm and  $\alpha \in ]0, +\infty[, \beta \in ]0, +\infty[$ . Note that the purpose of the norm  $\|.\|_{(\alpha,\beta)}$  is to tag the contributions of perturbations of A and b in the condition number; see [3].

The purpose of this paper is to exhibit the square of the condition number of  $x_r$  as the 2-norm of a symmetric nonnegative matrix  $\Delta$  that can be formed from the singular values of A, and the Fourier coefficients given by the entries of  $U^T b$ . The paper is organized as follows. In section 2, we state preliminary results based on results from [12]. Section 3 is devoted to an expression for the first-order expansion of  $x_r$  with respect to the data (A, b). The main result of this section is the matrix representation for the corresponding Fréchet derivative leading to the formula for the condition number of  $x_r$  using the singular values of A and the Fourier coefficients of b. We perform some numerical tests to validate our analysis by comparing it with results produced by a finite difference approach in section 4. A brief conclusion is given in section 5.

**2.** Preliminary results. It will be worthwhile to define the following matrix partitions:

$$V = [V_r, V_{\perp}] \in \Re^{p \times p}, \quad U = [U_r, U_{\perp}] \in \Re^{n \times n}, \quad \Sigma = \begin{bmatrix} \Sigma_r \\ \Sigma_{\perp} \end{bmatrix} \in \Re^{n \times p},$$

where

$$V_r \in \Re^{p \times r}, \qquad V_{\perp} \in \Re^{p \times (p-r)}, \qquad U_r \in \Re^{n \times r}, \qquad U_{\perp} \in \Re^{n \times (n-r)},$$
$$\Sigma_r = \operatorname{diag}(\sigma_1, \dots, \sigma_r) \in \Re^{r \times r}, \quad \Sigma_{\perp} = \begin{bmatrix} \operatorname{diag}(\sigma_{r+1}, \dots, \sigma_p) \\ 0 \end{bmatrix} \in \Re^{(n-r) \times (p-r)},$$

Furthermore, we define matrices  $E_{rr} = U_r^T E V_r$ ,  $E_{r\perp} = U_r^T E V_{\perp}$ ,  $E_{\perp r} = U_{\perp}^T E V_r$ , and  $E_{\perp\perp} = U_{\perp}^T E V_{\perp}$  and vectors  $b_r = U_r^T b$ ,  $b_{\perp} = U_{\perp}^T b$ , and  $f_r = U_r^T f$ . Finally, we shall denote by  $I_r$ ,  $I_{n-r}$  and  $I_{p-r}$  the identity matrices of order r, n-r, and p-r, respectively.

The operator vec (·) and the Kronecker product  $\otimes$  will be of particular importance in what follows. The vec (·) operator stacks the columns of the matrix argument into one long vector. For any matrices B and C, the matrix  $B \otimes C = (b_{ij}C)$ . It is enough for our purpose to recall the following properties concerning these operators.<sup>1</sup> For any matrices B, X, and C having compatible dimensions with respect to the involved products, we have

(2.1) 
$$\operatorname{vec}(BXC) = (C^T \otimes B)\operatorname{vec}(X),$$

(2.2) 
$$\operatorname{vec}(X^T) = \Psi_{(n,p)}\operatorname{vec}(X) \text{ for all } X \in \mathfrak{R}^{n \times p},$$

where  $\Psi_{(n,p)} \in \Re^{np \times np}$  is the permutation matrix defined by

$$\Psi_{(n,p)} = \sum_{i=1}^{n} \sum_{j=1}^{p} L_{ij} \otimes L_{ij}^{T}$$

Here each  $L_{ij} \in \Re^{n \times p}$  has entry 1 in position (i, j) and all other entries are zero.

The following assumption will be of particular importance in what follows.

Assumption 2.1. Let

$$\gamma = \| (E_{\perp r}^T, E_{r\perp}) \|_F,$$

suppose that

$$\delta = |\sigma_r - \sigma_{r+1}| - ||E_{rr}||_2 - ||E_{\perp\perp}||_2 > 0,$$

and assume that

$$\gamma/\delta < 1/2.$$

Roughly speaking, the statement of Assumption 2.1 is that the existence of a gap between  $\sigma_r$  and  $\sigma_{r+1} > 0$  is required and that  $||E||_2$  must be small enough compared to this gap.

<sup>&</sup>lt;sup>1</sup>We refer to [10, Chapter 4] for further properties of these operators.

Now, we state and adapt results from [12] to our context in the following two theorems.

THEOREM 2.2 (see [12, Theorem 6.4]). Let an  $n \times p$  perturbation matrix E be given and partition  $U^T E V$  with respect to  $U = [U_r, U_\perp]$  and  $V = [V_r, V_\perp]$  in the form

$$U^T E V = \begin{pmatrix} E_{rr} & E_{r\perp} \\ E_{\perp r} & E_{\perp \perp} \end{pmatrix}$$

Then under Assumption 2.1, there are matrices  $Q \in \Re^{(n-r) \times r}$  and  $P \in \Re^{(p-r) \times r}$  satisfying

$$\|(Q^T, P^T)\|_F < 2\frac{\gamma}{\delta} < 1$$

such that range $(V_r + V_{\perp}P)$  and range $(U_r - U_{\perp}Q)$  form a pair of singular subspaces for  $\tilde{A} = A + E$ .

Among other things, the theorem above tells us that Q and P approach 0 as E approaches 0. Other useful results related to the ones above are given in the following theorem (see again [12] and [13, p. 266]).

THEOREM 2.3. Suppose Assumption 2.1 holds. Then there exist matrices  $Q \in \Re^{(n-r) \times r}$  and  $P \in \Re^{(p-r) \times r}$  such that

(2.3) 
$$\tilde{U}_r = (U_r - U_\perp Q)(I + Q^T Q)^{-1/2}, \quad \tilde{U}_\perp = (U_r Q^T + U_\perp)(I + Q Q^T)^{-1/2},$$
  
(2.4)  $\tilde{V}_r = (V_r + V_\perp P)(I + P^T P)^{-1/2}, \quad \tilde{V}_\perp = (-V_r P^T + V_\perp)(I + P P^T)^{-1/2},$ 

with  $\tilde{U}_r^T \tilde{A} \tilde{V}_{\perp} = 0$  and  $\tilde{U}_{\perp}^T \tilde{A} \tilde{V}_r = 0$ . Furthermore,  $\tilde{U} = [\tilde{U}_r, \tilde{U}_{\perp}] \in \Re^{n \times n}$  and  $\tilde{V} = [\tilde{V}_r, \tilde{V}_{\perp}] \in \Re^{p \times p}$  are orthogonal matrices.

Since the overall aim of this investigation is to derive the condition number as the norm of the Fréchet derivative of  $x_r$ , our intermediate goal will be to write a first-order expansion of (1.1) in terms of quantities in (2.3) and (2.4) and then replace Q and P with their respective first-order expansions with respect to E. The next theorem exploits (2.3) and (2.4) together with properties of singular decomposition to establish these expansions.

THEOREM 2.4. Suppose that  $\sigma_r - \sigma_{r+1} > 0$ . Then the first-order expansions for Q and P are given by

(2.5) 
$$\operatorname{vec} (\mathbf{Q}^{\mathrm{T}}) = -\left(I_{n-r} \otimes \Sigma_{r}^{2} - (\Sigma_{\perp} \Sigma_{\perp}^{T}) \otimes I_{r}\right)^{-1} \times \left[I_{n-r} \otimes \Sigma_{r}, \Sigma_{\perp} \otimes I_{r}\right] \begin{bmatrix} \Psi_{(n-r,r)}(V_{r}^{T} \otimes U_{\perp}^{T}) \\ V_{\perp}^{T} \otimes U_{r}^{T} \end{bmatrix} \operatorname{vec}(\mathbf{E}) + \operatorname{o}(||E||),$$
  
(2.6) 
$$\operatorname{vec} (\mathbf{P}) = \left(\Sigma_{r}^{2} \otimes I_{p-r} - I_{r} \otimes (\Sigma_{\perp}^{T} \Sigma_{\perp})\right)^{-1} \times \left[I_{r} \otimes \Sigma_{\perp}^{T}, \Sigma_{r} \otimes I_{p-r}\right] \begin{bmatrix} V_{r}^{T} \otimes U_{\perp}^{T} \\ \Psi_{(r,p-r)} (V_{\perp}^{T} \otimes U_{r}^{T}) \end{bmatrix} \operatorname{vec}(\mathbf{E}) + \operatorname{o}(||E||).$$

*Proof.* In agreement with

(2.7) 
$$U^T A V = \begin{bmatrix} U_r^T A V_r & U_r^T A V_\perp \\ U_\perp^T A V_r & U_\perp^T A V_\perp \end{bmatrix} = \begin{bmatrix} \Sigma_r & 0 \\ 0 & \Sigma_\perp \end{bmatrix} \in \Re^{n \times p},$$

together with the results of Theorem 2.3, we have

(2.8) 
$$U^{T}\tilde{A}V = \begin{bmatrix} U_{r}^{T}(A+E)V_{r} & U_{r}^{T}(A+E)V_{\perp} \\ U_{\perp}^{T}(A+E)V_{r} & U_{\perp}^{T}(A+E)V_{\perp} \end{bmatrix}$$
$$\stackrel{\text{def}}{=} \begin{bmatrix} \Sigma_{r}+E_{rr} & E_{r\perp} \\ E_{\perp r} & \Sigma_{\perp}+E_{\perp\perp} \end{bmatrix},$$
$$\begin{bmatrix} \widetilde{\omega}T, \widetilde{\omega}, -\widetilde{\omega}T, -$$

(2.9) 
$$\tilde{U}^T \tilde{A} \tilde{V} = \begin{bmatrix} \tilde{U}_r^T \tilde{A} \tilde{V}_r & \tilde{U}_r^T \tilde{A} \tilde{V}_\perp \\ \tilde{U}_\perp^T \tilde{A} \tilde{V}_r & \tilde{U}_\perp^T \tilde{A} \tilde{V}_\perp \end{bmatrix} = \begin{bmatrix} \star & 0 \\ 0 & \star \end{bmatrix}$$

If we substitute (2.3)–(2.4) into the extra diagonal blocks of (2.9) (that are zero), we obtain

.

$$\begin{aligned} &-\left(QU_r^TAV_r + QU_r^TAV_{\perp}P + QU_r^TEV_r + QU_r^TEV_{\perp}P\right) \\ &-U_{\perp}^TAV_r - U_{\perp}^TAV_{\perp}P - U_{\perp}^TEV_r - U_{\perp}^TEV_{\perp}P) = 0, \\ &-\left(U_r^TAV_rP^T - U_r^TAV_{\perp} + U_r^TEV_rP^T - U_r^TEV_{\perp}\right) \\ &+ U_{\perp}^TAV_rP^T - Q^TU_{\perp}^TAV_{\perp} + Q^TU_{\perp}^TEV_rP^T - Q^TU_{\perp}^TEV_{\perp}) = 0. \end{aligned}$$

Furthermore, using relations (2.7) and (2.8) and after rearranging terms, we obtain (see also [12, equation (6.2)]) the pair of quadratic matrix equations

(2.10) 
$$Q(\Sigma_r + E_{rr}) + (\Sigma_\perp + E_{\perp\perp})P = -E_{\perp r} - QE_{r\perp}P,$$

(2.11) 
$$P(\Sigma_r + E_{rr}^T) + (\Sigma_{\perp}^T + E_{\perp\perp}^T)Q = E_{r\perp}^T + PE_{\perp r}^TQ,$$

where unknowns are Q and P. We retain only first-order terms  $^2$  in  $\|E\|$  in (2.10) and (2.11) leading to

(2.12) 
$$Q\Sigma_r + \Sigma_\perp P = -E_{\perp r} + o(||E||),$$

(2.13) 
$$P\Sigma_r + \Sigma_{\perp}^T Q = E_{r\perp}^T + o(||E||),$$

from which we obtain the system

(2.14) 
$$Q = -\Sigma_{\perp} P \Sigma_r^{-1} - E_{\perp r} \Sigma_r^{-1} + o(||E||),$$

(2.15) 
$$P = -\Sigma_{\perp}^{T} Q \Sigma_{r}^{-1} + E_{r \perp}^{T} \Sigma_{r}^{-1} + o(||E||)$$

by a postmultiplication of both (2.12) and (2.13) by  $\Sigma_r$  (which exists because  $\sigma_1 \geq \cdots \geq \sigma_r > \sigma_{r+1} \geq 0$ ). Replacing P in (2.14) by the right-hand side of (2.15) and conversely replacing Q in (2.15) by the right-hand side of (2.14) we have

(2.16) 
$$Q = -\Sigma_{\perp} (-\Sigma_{\perp}^T Q \Sigma_r^{-1} + E_{r\perp}^T \Sigma_r^{-1}) \Sigma_r^{-1} - E_{\perp r} \Sigma_r^{-1} + o(||E||),$$

(2.17) 
$$P = -\Sigma_{\perp}^{T} (-\Sigma_{\perp} P \Sigma_{r}^{-1} - E_{\perp r} \Sigma_{r}^{-1}) \Sigma_{r}^{-1} + E_{r\perp}^{T} \Sigma_{r}^{-1} + o(||E||).$$

Postmultiplying (2.16) and (2.17) by  $\Sigma_r^2$  and rearranging terms yields

(2.18) 
$$\Sigma_r^2 Q^T - Q^T \Sigma_{\perp} \Sigma_{\perp}^T = -E_{r\perp} \Sigma_{\perp}^T - \Sigma_r E_{\perp r}^T + o(||E||),$$

(2.19) 
$$P\Sigma_r^2 - \Sigma_{\perp}^T \Sigma_{\perp} P = \Sigma_{\perp}^T E_{\perp r} + E_{r\perp}^T \Sigma_r + o(||E||).$$

<sup>2</sup>This is why the terms  $PE_{rr}^T$ ,  $E_{\perp\perp}^T Q$ ,  $PE_{\perp r}^T Q$ ,  $QE_{rr}$ ,  $E_{\perp\perp}P$ , and  $QE_{r\perp}P$  no longer appear.

According to property (2.1), (2.18) and (2.19) may be rewritten as

$$\begin{aligned} & \left(I_{n-r} \otimes \Sigma_r^2 - (\Sigma_{\perp} \Sigma_{\perp}^T) \otimes I_r\right) \operatorname{vec}(Q^T) \\ &= -\operatorname{vec}\left(E_{r\perp} \Sigma_{\perp}^T + \Sigma_r E_{\perp r}^T\right) + \operatorname{o}(\|E\|) \\ &= -\left[I_{n-r} \otimes \Sigma_r, \Sigma_{\perp} \otimes I_r\right] \begin{bmatrix} \operatorname{vec}\left(E_{\perp r}^T\right) \\ \operatorname{vec}\left(E_{r\perp}\right) \end{bmatrix} + \operatorname{o}(\|E\|), \\ & \left(\Sigma_r^2 \otimes I_{p-r} - I_r \otimes (\Sigma_{\perp}^T \Sigma_{\perp})\right) \operatorname{vec}(P) \\ &= \operatorname{vec}\left(\Sigma_{\perp}^T E_{\perp r} + E_{r\perp}^T \Sigma_r\right) + \operatorname{o}(\|E\|) \\ &= \left[I_r \otimes \Sigma_{\perp}^T, \Sigma_r \otimes I_{p-r}\right] \begin{bmatrix} \operatorname{vec}\left(E_{\perp r}\right) \\ \operatorname{vec}\left(E_{r\perp}^T\right) \end{bmatrix} + \operatorname{o}(\|E\|). \end{aligned}$$

One can replace  $\operatorname{vec}(E_{\perp r}^T)$  and  $\operatorname{vec}(E_{r\perp}^T)$  by  $\Psi(n-r,r)\operatorname{vec}(E_{\perp r})$  and  $\Psi(r,p-r)\operatorname{vec}(E_{r\perp})$ , respectively, based on property (2.2). Note that  $(I_{n-r} \otimes \Sigma_r^2 - (\Sigma_{\perp} \Sigma_{\perp}^T) \otimes I_r)$  and  $(\Sigma_r^2 \otimes I_{p-r} - I_r \otimes (\Sigma_{\perp}^T \Sigma_{\perp}))$  are diagonal matrices of order (n-r)r and (p-r)r, respectively. In addition, their diagonal entries are strictly positive since  $\sigma_r > \sigma_{r+1}$ . Hence, their inverses exist. To conclude the proof, observe that

$$\operatorname{vec}(E_{\perp r}) = (V_r^T \otimes U_{\perp}^T) \operatorname{vec}(E), \qquad \operatorname{vec}(E_{r\perp}) = (V_{\perp}^T \otimes U_r^T) \operatorname{vec}(E),$$
$$\operatorname{vec}(E_{\perp\perp}) = (V_{\perp}^T \otimes U_{\perp}^T) \operatorname{vec}(E), \qquad \operatorname{vec}(E_{rr}) = (V_r^T \otimes U_r^T) \operatorname{vec}(E). \qquad \Box$$

In what follows, we use the results in Theorem 2.3 to introduce the first-order expansion for  $x_r$  around (A, b) in terms of the partitioned singular value decomposition matrices of A, the perturbation matrix E, the vector b, and the perturbation vector f.

3. The Fréchet derivative and the condition number of  $x_r$ . The continuity and the differentiability of  $x_r$  rely on the fact that one supposes that there is a gap between  $\sigma_r$  and  $\sigma_{r+1}$ , that is,  $\sigma_r - \sigma_{r+1} > 0$ . Consider the following counterexample. Let

$$A = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad E = \begin{pmatrix} \epsilon^2 \sin(\frac{1}{\epsilon}) & 0 \\ 0 & -\epsilon^2 \sin(\frac{1}{\epsilon}) \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad f = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$

We take r = 1. Thus

$$\tilde{x}_r = \begin{cases} \frac{1}{1+\epsilon^2 \sin(\frac{1}{\epsilon})} e_1 & \text{if } \sin(\frac{1}{\epsilon}) > 0, \\ \frac{1}{1-\epsilon^2 \sin(\frac{1}{\epsilon})} e_2 & \text{if } \sin(\frac{1}{\epsilon}) < 0, \\ x_r & \text{if } \sin(\frac{1}{\epsilon}) = 0, \end{cases}$$

where  $e_1 = (1,0)^T$  and  $e_2 = (0,1)^T$  are the canonical vectors of  $\Re^2$ . The above counterexample shows that the unit-vector of  $\tilde{x}_r$  fluctuates between  $e_1$  and  $e_2$  as  $\epsilon$ tends to 0. In this case  $x_r$  is not continuous, and a fortiori not differentiable, around A. We know from Theorem 2.3 that the singular values of  $\tilde{A}$  are the disjoint union of the singular values of  $\tilde{U}_r^T \tilde{A} \tilde{V}_r$  and those of  $\tilde{U}_\perp^T \tilde{A} \tilde{V}_\perp$ . To define  $\tilde{x}_r$  by (1.1) it is required that the *r* leading singular values of  $\tilde{A}$  be those of  $\tilde{U}_r^T \tilde{A} \tilde{V}_r$ . This is achieved if  $\sigma_r - \sigma_{r+1} > 0$ and E, sufficiently small.<sup>3</sup>

<sup>&</sup>lt;sup>3</sup>Observe that in the presence of a gap  $\sigma_r - \sigma_{r+1} > 0$ , the bases of the involved singular subspaces of  $\tilde{A}$  tend continuously to those of A as E tends 0.

Now, let us state the following lemma.

LEMMA 3.1. Suppose  $\sigma_r - \sigma_{r+1} > 0$ . Then the first-order expansion of  $x_r$  can be written in the form

(3.1) 
$$\tilde{x}_r = x_r + V \begin{bmatrix} I_r \\ 0 \end{bmatrix} \Sigma_r^{-1} f_r - V \begin{bmatrix} I_r \\ 0 \end{bmatrix} \Sigma_r^{-1} Q^T b_\perp + V \begin{bmatrix} 0 \\ I_{p-r} \end{bmatrix} P \Sigma_r^{-1} b_r - V \begin{bmatrix} I_r \\ 0 \end{bmatrix} \Sigma_r^{-1} E_{rr} \Sigma_r^{-1} b_r + o(||[E, f]||).$$

*Proof.* We insert (2.3) and (2.4) in expression (1.1) to yield

$$\tilde{x}_r = (V_r + V_{\perp}P)((U_r - U_{\perp}Q)^T (A + E)(V_r + V_{\perp}P))^{-1}(U_r - U_{\perp}Q)^T \tilde{b} = (V_r + V_{\perp}P)(\Sigma_r^{-1} - \Sigma_r^{-1}U_r^T E V_r \Sigma_r^{-1})(U_r - U_{\perp}Q)^T \tilde{b} + o(\|[E, f]\|),$$

where we used the following result concerning a perturbation of the inverse of a matrix  $(F+G)^{-1} = F^{-1} - F^{-1}GF^{-1} + o(||G||)$ ; see [13, p. 131]. Developing this equation and recalling that  $E_{rr} \stackrel{\text{def}}{=} U_r^T EV_r$  gives, after rearranging terms,

$$\begin{split} \tilde{x}_r &= x_r + V_r \Sigma_r^{-1} U_r^T f - V_r \Sigma_r^{-1} Q^T U_{\perp}^T b + V_{\perp} P \Sigma_r^{-1} U_r^T b - V_r \Sigma_r^{-1} E_{rr} \Sigma_r^{-1} U_r^T b \\ &+ \mathrm{o}(\|[E, f]\|) \\ &= x_r + V_r \Sigma_r^{-1} f_r - V_r \Sigma_r^{-1} Q^T b_{\perp} + V_{\perp} P \Sigma_r^{-1} b_r - V_r \Sigma_r^{-1} E_{rr} \Sigma_r^{-1} b_r + \mathrm{o}(\|[E, f]\|) \end{split}$$

From the properties

$$VV^T = I, V^T V_r = \begin{bmatrix} I_r \\ 0 \end{bmatrix}$$
 and  $V^T V_\perp = \begin{bmatrix} 0 \\ I_{p-r} \end{bmatrix},$ 

we have

$$\tilde{x}_r = x_r + VV^T V_r \Sigma_r^{-1} f_r - VV^T V_r \Sigma_r^{-1} Q^T b_\perp + VV^T V_\perp P \Sigma_r^{-1} b_r - VV^T V_r \Sigma_r^{-1} E_{rr} \Sigma_r^{-1} b_r + o(||[E, f]||),$$

which implies (3.1).

Now, we are ready to give the expression of the matrix  $x'_r$  that represents the Fréchet derivative of  $x_r$ , with respect to the data (A, b). The expression is given in terms of the singular value decomposition information of A and the vector b. For that, we simply use results in Theorem 2.4 to eliminate Q and P from (3.1).

PROPOSITION 3.2. Suppose that  $\sigma_r - \sigma_{r+1} > 0$ . Then the application

$$x_r : (\Re^{n \times p}, \Re^n) \longrightarrow \Re^p : (A, b) \longrightarrow x_r$$

is a differentiable function of (A, b). In addition, we have

$$\tilde{x}_r = x_r + x'_r \begin{bmatrix} \alpha \operatorname{vec}(\mathbf{E}) \\ \beta f \end{bmatrix} + \operatorname{o}(\|[E, f]\|)$$

with

(3.2) 
$$x'_r = V \begin{bmatrix} \frac{1}{\alpha} M, \frac{1}{\beta} \begin{pmatrix} \Sigma_r^{-1} \\ 0 \end{bmatrix} \end{bmatrix} W \in \Re^{n \times (np+n)}.$$

Here, W is an orthogonal matrix defined by

$$W = \begin{bmatrix} V_r^T \otimes U_{\perp}^T & 0 \\ V_{\perp}^T \otimes U_r^T & 0 \\ V_r^T \otimes U_r^T & 0 \\ V_{\perp}^T \otimes U_{\perp}^T & 0 \\ 0 & U^T \end{bmatrix} \in \Re^{(np+n) \times (np+n)},$$

and M is the partitioned matrix given by

$$M = \begin{bmatrix} R_r & S_r & -T_r & 0\\ R_\perp & S_\perp & 0 & 0 \end{bmatrix} \in \Re^{p \times (np)}$$

with

$$(3.3) \quad R_r = (b_{\perp}^T \otimes \Sigma_r^{-1}) \left( I_{n-r} \otimes \Sigma_r^2 - (\Sigma_{\perp} \Sigma_{\perp}^T) \otimes I_r \right)^{-1} (I_{n-r} \otimes \Sigma_r) \Psi_{(n-r,r)},$$

$$(3.4) \quad S_r = (b_{\perp}^T \otimes \Sigma_r^{-1}) \left( I_{n-r} \otimes \Sigma_r^2 - (\Sigma_{\perp} \Sigma_{\perp}^T) \otimes I_r \right)^{-1} (\Sigma_{\perp} \otimes I_r),$$

$$(3.5) \quad R_{\perp} = \left( (b_r^T \Sigma_r^{-1}) \otimes I_{p-r} \right) \left( \Sigma_r^2 \otimes I_{p-r} - I_r \otimes (\Sigma_{\perp}^T \Sigma_{\perp}) \right)^{-1} (I_r \otimes \Sigma_{\perp}^T),$$

$$(3.6) \quad S_{\perp} = \left( (b_r^T \Sigma_r^{-1}) \otimes I_{p-r} \right) \left( \Sigma_r^2 \otimes I_{p-r} - I_r \otimes (\Sigma_{\perp}^T \Sigma_{\perp}) \right)^{-1} (\Sigma_r \otimes I_{p-r}) \Psi_{(r,p-r)},$$

$$(3.7) \quad T_r = \left( b_r^T \Sigma_r^{-1} \right) \otimes \Sigma_r^{-1}.$$

The dimensions of these matrices are given in the following:

$$R_r, S_r \in \Re^{r \times (n-r)r}, \qquad R_\perp, S_\perp \in \Re^{(p-r) \times (n-r)r}, \qquad and T_r \in \Re^{r \times r^2}.$$

*Proof.* Consider the quantities in (3.1). Using the properties of the vec operator applied to a vector, we obtain

$$\begin{bmatrix} I_r \\ 0 \end{bmatrix} \Sigma_r^{-1} E_{rr} \Sigma_r^{-1} b_r = \begin{bmatrix} (b_r^T \Sigma_r^{-1}) \otimes \Sigma_r^{-1} \\ 0 \end{bmatrix} \operatorname{vec} (E_{rr}) = \begin{bmatrix} T_r \\ 0 \end{bmatrix} (V_r^T \otimes U_r^T) \operatorname{vec} (E).$$

Taking the expressions for  $vec(Q^T)$  and vec(P) given in (2.5) and (2.6), we have

$$\begin{bmatrix} I_r \\ 0 \end{bmatrix} \Sigma_r^{-1} Q^T b_\perp = \begin{bmatrix} b_\perp^T \otimes \Sigma_r^{-1} \\ 0 \end{bmatrix} \operatorname{vec} (Q^T)$$
$$= -\begin{bmatrix} R_r & S_r \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_r^T \otimes U_\perp^T \\ V_\perp^T \otimes U_r^T \end{bmatrix} \operatorname{vec} (E) + o(\|[E, f]\|),$$
$$\begin{bmatrix} 0 \\ I_{p-r} \end{bmatrix} P \Sigma_r^{-1} b_r = \begin{bmatrix} 0 \\ (b_r^T \Sigma_r^{-1}) \otimes I_{p-r} \end{bmatrix} \operatorname{vec} (P)$$
$$= \begin{bmatrix} 0 & 0 \\ R_\perp & S_\perp \end{bmatrix} \begin{bmatrix} V_r^T \otimes U_\perp^T \\ V_\perp^T \otimes U_r^T \end{bmatrix} \operatorname{vec} (E) + o(\|[E, f]\|).$$

Injecting these quantities in (3.1) results in

$$\tilde{x}_{r} = x_{r} + V \begin{bmatrix} \Sigma_{r}^{-1} \\ 0 \end{bmatrix} U^{T} f + V \begin{bmatrix} R_{r} & S_{r} & -T_{r} & 0 \\ R_{\perp} & S_{\perp} & 0 & 0 \end{bmatrix} \begin{bmatrix} V_{r}^{T} \otimes U_{\perp}^{T} \\ V_{\perp}^{T} \otimes U_{r}^{T} \\ V_{r}^{T} \otimes U_{r}^{T} \\ V_{\perp}^{T} \otimes U_{\perp}^{T} \end{bmatrix}$$
vec (E)

+ o(||[E, f]||),

from which the results are derived.  $\hfill \Box$ 

We can now establish the expression of the  $x_r$  condition number. We know by definition that

$$|||x'_r|||_{(\alpha,\beta),2} = \max_{[\alpha E,\beta f]\neq 0} \frac{||x'_r \cdot (E,f)||_2}{||\operatorname{vec} [E,f]||_{(\alpha,\beta)}}.$$

Thus, from (3.2) we conclude that the exact condition number of  $x_r$  is

$$|||x'_r|||_{(\alpha,\beta),2} = \lambda_{\max}^{1/2}(\Delta),$$

where

$$\Delta \stackrel{\text{def}}{=} V^T x'_r (x'_r)^T V = \frac{1}{\alpha^2} M M^T + \frac{1}{\beta^2} \begin{pmatrix} \Sigma_r^{-2} & 0\\ 0 & 0 \end{pmatrix} \in \Re^{p \times p}.$$

It remains to show how  $\Delta$  can be expressed with the singular values of A and the Fourier coefficients given by the elements of  $U^T b$ .

PROPOSITION 3.3. Assume that the singular values of the matrix A are such that

$$\sigma_1 \ge \cdots \ge \sigma_r > \sigma_{r+1} \ge \cdots \ge \sigma_p \ge 0.$$

Then

$$\Delta = \begin{pmatrix} \frac{1}{\alpha^2} \Delta_{rr} + \frac{1}{\beta^2} \Sigma_r^{-2} & \frac{1}{\alpha^2} \Gamma_{\perp r}^T \\ \frac{1}{\alpha^2} \Gamma_{\perp r} & \frac{1}{\alpha^2} \Delta_{\perp \perp} \end{pmatrix},$$

where

$$\begin{split} \Delta_{rr} &= \operatorname{diag}\left(\sum_{k=1}^{r} \frac{\theta_k^2}{\sigma_k^2 \sigma_t^2} + \sum_{k=r+1}^{p} (\pi_k^{(t)})^2 \frac{\sigma_k^2 + \sigma_t^2}{\sigma_t^2} \theta_k^2 + \sum_{k=p+1}^{n} \frac{\theta_k^2}{\sigma_t^4}\right), \quad 1 \le t \le r, \\ \Delta_{\perp\perp} &= \operatorname{diag}\left(\sum_{k=1}^{r} (\pi_k^{(t)})^2 \frac{\sigma_k^2 + \sigma_t^2}{\sigma_k^2} \theta_k^2\right), \qquad r+1 \le t \le p, \end{split}$$

$$\begin{split} \Gamma_{\perp r} &= R_{\perp} R_{r}^{T} + S_{\perp} S_{r}^{T} \\ &= 2 \begin{pmatrix} (\pi_{r+1}^{(1)})^{2} \frac{\sigma_{r+1}}{\sigma_{1}} \theta_{1} \theta_{r+1} & (\pi_{r+1}^{(2)})^{2} \frac{\sigma_{r+1}}{\sigma_{2}} \theta_{2} \theta_{r+1} & \cdots & (\pi_{r+1}^{(r)})^{2} \frac{\sigma_{r+2}}{\sigma_{r}} \theta_{r} \theta_{r+1} \\ (\pi_{r+2}^{(1)})^{2} \frac{\sigma_{r+2}}{\sigma_{1}} \theta_{1} \theta_{r+2} & (\pi_{r+2}^{(2)})^{2} \frac{\sigma_{r+2}}{\sigma_{2}} \theta_{2} \theta_{r+2} & \cdots & (\pi_{r+2}^{(r)})^{2} \frac{\sigma_{r+2}}{\sigma_{r}} \theta_{r} \theta_{r+2} \\ &\vdots & \vdots & \ddots & \vdots \\ (\pi_{p}^{(1)})^{2} \frac{\sigma_{p}}{\sigma_{1}} \theta_{1} \theta_{p} & (\pi_{p}^{(2)})^{2} \frac{\sigma_{p}}{\sigma_{2}} \theta_{2} \theta_{p} & \cdots & (\pi_{p}^{(r)})^{2} \frac{\sigma_{p}}{\sigma_{r}} \theta_{r} \theta_{p} \end{pmatrix}, \end{split}$$

with  $(\theta_1, \ldots, \theta_n) = b^T U$ , and  $\pi_k^{(t)} = \frac{1}{\sigma_t^2 - \sigma_k^2}$ , with either  $t = 1, \ldots, r$  and  $k = r+1, \ldots, p$ or  $k = 1, \ldots, r$  and  $t = r+1, \ldots, p$ .

Moreover, the quantity  $\pi_k^{(t)}$  is well defined, since whenever it appears,  $\sigma_t^2 - \sigma_k^2 \neq 0$  holds.

*Proof.* First we consider the  $p \times p$  symmetric matrix

$$MM^{T} = \begin{bmatrix} R_{r}R_{r}^{T} + S_{r}S_{r}^{T} + T_{r}T_{r}^{T} & -R_{r}R_{\perp}^{T} - S_{r}S_{\perp}^{T} \\ -R_{\perp}R_{r}^{T} - S_{\perp}S_{r}^{T} & R_{\perp}R_{\perp}^{T} + S_{\perp}S_{\perp}^{T} \end{bmatrix} \stackrel{\text{def}}{=} \begin{bmatrix} \Delta_{rr} & \Gamma_{\perp r} \\ \Gamma_{r\perp} & \Delta_{\perp\perp} \end{bmatrix}.$$

Exploiting their structure, we can write the matrices (3.3)–(3.7) as

 $\begin{array}{ll} (3.8) & R_r = \left[\theta_{r+1}(\varSigma_r^2 - \sigma_{r+1}^2 I_r)^{-1}, \dots, \theta_p(\varSigma_r^2 \\ (3.9) & -\sigma_p^2 I_r)^{-1}, \theta_{p+1} \varSigma_r^{-2}, \dots, \theta_n \varSigma_r^{-2}\right] \Psi_{(n-r,r)}, \\ (3.10) & S_r = \left[\theta_{r+1} \sigma_{r+1} \varSigma_r^{-1}(\varSigma_r^2 - \sigma_{r+1}^2 I_r)^{-1}, \dots, \theta_p \sigma_p \varSigma_r^{-1}(\varSigma_r^2 - \sigma_p^2 I_r)^{-1}, 0, \dots, 0\right], \\ (3.11) & R_\perp = \left[\theta_1 \sigma_1^{-1} (\sigma_1^2 I_{p-r} - \varSigma_\perp^T \varSigma_\perp)^{-1} \varSigma_\perp^T, \dots, \theta_r \sigma_r^{-1} (\sigma_r^2 I_{p-r} - \varSigma_\perp^T \varSigma_\perp)^{-1} \varSigma_\perp^T\right], \\ (3.12) & S_\perp = \left[\theta_1 (\sigma_1^2 I_{p-r} - \varSigma_\perp^T \varSigma_\perp)^{-1}, \dots, \theta_r (\sigma_r^2 I_{p-r} - \varSigma_\perp^T \varSigma_\perp)^{-1}\right] \Psi_{(r,p-r)}, \\ (3.13) & T_r = \left[\theta_1 \sigma_1^{-1} \varSigma_r^{-1}, \dots, \theta_r \sigma_r^{-1} \varSigma_r^{-1}\right]. \end{array}$ 

In (3.8), the first of the two factors,

(3.14) 
$$\left[ \theta_{r+1} (\Sigma_r^2 - \sigma_{r+1}^2 I_r)^{-1}, \dots, \theta_p (\Sigma_r^2 - \sigma_p^2 I_r)^{-1}, \theta_{p+1} \Sigma_r^{-2}, \dots, \theta_n \Sigma_r^{-2} \right]$$

is a  $1\times(n-r)$  partitioned matrix. Its blocks consist of  $r\text{-}\mathrm{order}$  diagonal matrices. Recall that the second factor in (3.8) is

(3.15) 
$$\Psi_{(n-r,r)} = \sum_{i=1}^{n-r} \sum_{j=1}^{r} L_{ij} \otimes L_{ij}^{T},$$

where  $L_{ij} = e_i e_j^T \in \Re^{(n-r)r}$  with  $e_i \in \Re^{(n-r)}$  and  $e_j \in \Re^r$ . Observe that  $L_{ij} \otimes L_{ij}^T$  is an  $(n-r) \times r$  partitioned matrix where each block has r rows and n-r columns. Furthermore, it has the block  $L_{ij}^T$  in position i, j and 0 in the remaining blocks. The multiplication of the partitioned matrices (3.14) and (3.15) results in the  $1 \times r$  partitioned matrix

$$R_{r} = \sum_{i=1}^{n-r} \sum_{j=1}^{r} \left[ \theta_{r+1} (\Sigma_{r}^{2} - \sigma_{r+1}^{2} I_{r})^{-1}, \dots, \theta_{p} (\Sigma_{r}^{2} - \sigma_{p}^{2} I_{r})^{-1}, \theta_{p+1} \Sigma_{r}^{-2}, \dots, \theta_{n} \Sigma_{r}^{-2} \right] L_{ij} \otimes L_{ij}^{T},$$

whose block j can be written as

$$\sum_{i=1}^{p-r} \theta_{r+i} (\Sigma_r^2 - \sigma_{r+i}^2 I_r)^{-1} L_{ij}^T + \sum_{i=p-r+1}^{n-r} \theta_{r+i} \Sigma_r^{-2} L_{ij}^T$$

Consequently, multiplying  $R_{\perp}$  and  $R_r$  block by block yields

(3.16) 
$$R_{\perp}R_{r}^{T} = \sum_{j=1}^{r} \theta_{j}\sigma_{j}^{-1}(\sigma_{j}^{2}I_{p-r} - \Sigma_{\perp}^{T}\Sigma_{\perp})^{-1}\Sigma_{\perp}^{T}\sum_{i=1}^{p-r}L_{ij}\theta_{r+i}(\Sigma_{r}^{2} - \sigma_{r+i}I_{r})^{-1} + \sum_{j=1}^{r}\theta_{j}\sigma_{j}^{-1}(\sigma_{j}^{2}I_{p-r} - \Sigma_{\perp}^{T}\Sigma_{\perp})^{-1}\Sigma_{\perp}^{T}\sum_{i=p-r+1}^{n-r}L_{ij}\theta_{r+i}\Sigma_{r}^{-2}.$$

Since  $\Sigma_{\perp}^T \sum_{i=p-r+1}^{n-r} e_i = 0$ , one has  $\Sigma_{\perp}^T \sum_{i=p-r+1}^{n-r} L_{ij} = \Sigma_{\perp}^T \sum_{i=p-r+1}^{n-r} e_i e_j^T = 0$  and hence the last term in (3.16) vanishes. Thus

$$R_{\perp}R_{r}^{T} = \sum_{i=1}^{p-r} \sum_{j=1}^{r} \theta_{r+i}\theta_{j}\sigma_{j}^{-1}(\sigma_{j}^{2}I_{p-r} - \Sigma_{\perp}^{T}\Sigma_{\perp})^{-1}\Sigma_{\perp}^{T}e_{i}e_{j}^{T}(\Sigma_{r}^{2} - \sigma_{r+i}^{2}I_{r})^{-1}.$$

A direct computation gives

$$\theta_{r+i}\theta_j\sigma_j^{-1}(\sigma_j^2 I_{p-r} - \Sigma_{\perp}^T \Sigma_{\perp})^{-1}\Sigma_{\perp}^T e_i = \begin{cases} \frac{\theta_{r+i}\theta_j\sigma_{r+i}}{\sigma_j(\sigma_j^2 - \sigma_{r+i}^2)}e_i, & i = 1,\dots, p-r, \\ 0, & i = p-r+1,\dots, n-r, \end{cases}$$

where  $e_i \in \Re^{(n-r)}$  in the left-hand side and  $e_i \in \Re^{(p-r)}$  on the right-hand side. Then from

$$e_{j}^{T}(\Sigma_{r}^{2} - \sigma_{r+i}^{2}I_{r})^{-1} = \frac{1}{(\sigma_{j}^{2} - \sigma_{r+i}^{2})}e_{j}^{T},$$

where  $e_j \in \Re^r$  on both sides, we deduce that

$$R_{\perp}R_{r}^{T} = \sum_{i=1}^{p-r} \sum_{j=1}^{r} \frac{1}{(\sigma_{j}^{2} - \sigma_{r+i}^{2})^{2}} \frac{\sigma_{r+i}}{\sigma_{j}} \theta_{r+i} \theta_{j} e_{i} e_{j}^{T},$$
  
$$= \sum_{i=1}^{p-r} \sum_{j=1}^{r} (\pi_{r+i}^{(j)})^{2} \frac{\sigma_{r+i}}{\sigma_{j}} \theta_{r+i} \theta_{j} e_{i} e_{j}^{T} \in \Re^{(p-r) \times r}$$

with  $\pi_{r+i}^{(j)} = \frac{1}{(\sigma_j^2 - \sigma_{r+i}^2)}$ . In the same manner we can compute and show that  $S_\perp S_r^T$  is equivalent to  $R_\perp R_r^T$ .

The remaining blocks in  $MM^T$  are computed by performing the block matrixmatrix multiplications. So,

$$\begin{aligned} R_{r}R_{r}^{T} &= \sum_{k=r+1}^{p} \theta_{k}^{2} (\Sigma_{r}^{2} - \sigma_{k}^{2}I_{r})^{-2} + \sum_{k=p+1}^{n} \theta_{k}^{2} \Sigma_{r}^{-4} = \operatorname{diag}\left(\sum_{k=r+1}^{p} (\pi_{k}^{(t)})^{2} \theta_{k}^{2} + \sum_{k=r+1}^{n} \frac{\theta_{k}^{2}}{\sigma_{t}^{4}}\right), \\ S_{r}S_{r}^{T} &= \sum_{k=r+1}^{p} \theta_{k}^{2} \sigma_{k}^{2} \Sigma_{r}^{-2} (\Sigma_{r}^{2} - \sigma_{k}^{2}I_{r})^{-2} \\ &= \operatorname{diag}\left(\sum_{k=r+1}^{p} (\pi_{k}^{(t)})^{2} \frac{\sigma_{k}^{2}}{\sigma_{t}^{2}} \theta_{k}^{2}\right), \\ T_{r}T_{r}^{T} &= \sum_{k=1}^{r} \frac{\theta_{k}^{2}}{\sigma_{k}^{2}} \Sigma_{r}^{-2} \\ &= \operatorname{diag}\left(\sum_{k=1}^{r} \frac{\theta_{k}^{2}}{\sigma_{k}^{2} \sigma_{t}^{2}}\right) \end{aligned}$$

for t = 1, ..., r;

$$R_{\perp}R_{\perp}^{T} = \sum_{k=1}^{r} \frac{\theta_{k}^{2}}{\sigma_{k}^{2}} \Sigma_{\perp}^{T} \Sigma_{\perp} (\sigma_{k}^{2} I_{p-r} - \Sigma_{\perp}^{T} \Sigma_{\perp})^{-2} = \operatorname{diag}\left(\sum_{k=1}^{r} (\pi_{k}^{(t)})^{2} \frac{\sigma_{t}^{2}}{\sigma_{k}^{2}} \theta_{k}^{2}\right)$$
$$S_{\perp}S_{\perp}^{T} = \sum_{k=1}^{r} \theta_{k}^{2} (\sigma_{k}^{2} I_{p-r} - \Sigma_{\perp}^{T} \Sigma_{\perp})^{-2} = \operatorname{diag}\left(\sum_{k=1}^{r} (\pi_{k}^{(t)})^{2} \theta_{k}^{2}\right)$$

for t = 1, ..., p - r.

Putting the above results together yields the result.

TABLE 1 The exact value of  $cond(x_r)$  using the expression in Proposition 3.3 versus the finite difference estimate value using jacobianest for 12 problems.

Problem	$\operatorname{cond}(x_r)$ from 3.3	Finite difference estimate value of $cond(x_n)$	n	p	r
baart	7 156e+3	7.087e+3	20	20	5
blur	$2.516e \pm 1$	$2.516e \pm 1$	16	16	6
derive	1.698e + 3	$1.698e \pm 3$	12	12	10
foxgood	$2.896e \pm 1$	$2.896e \pm 1$	20	20	2
heat	4.486e + 1	4.478e+1	12	12	10
i_laplace	1.448e+4	1.367e + 4	20	20	7
parallax	1.412e + 5	1.411e + 5	26	12	10
phillips	5.731e + 1	5.731e + 1	12	12	10
shaw	1.044e + 3	1.044e + 3	12	12	8
spikes	8.178e + 2	8.178e + 2	12	12	4
ursell	3.716e + 5	3.716e + 5	20	20	3
wing	3.429e + 6	3.010e + 6	20	20	5

Let us point out the fact that an early result in [3], when r = p, that is, when we do not perform truncation (i.e., we assume that A is a full rank matrix), is a particular case of the results above. In fact, in this case,  $\Delta$  becomes diagonal and simplifies to

$$\Delta_{rr} = \operatorname{diag}\left(\sum_{k=1}^{p} \frac{\theta_k^2}{\sigma_k^2 \sigma_t^2} + \sum_{k=p+1}^{n} \frac{\theta_k^2}{\sigma_t^4}\right) = \operatorname{diag}\left(\frac{1}{\sigma_t^2} \left(\sum_{k=1}^{p} \frac{\theta_k^2}{\sigma_k^2} + \sum_{k=p+1}^{n} \frac{\theta_k^2}{\sigma_t^2}\right)\right)$$

for t = 1, ..., p. This implies the result given in [3], that is,

$$\begin{split} ||x_r'|||_{(\alpha,\beta),2} &= \sqrt{\frac{1}{\alpha^2} \frac{1}{\sigma_{min}^2} \left[ \sum_{k=1}^p \left( \frac{\theta_k}{\sigma_k} \right)^2 + \frac{1}{\sigma_{min}^2} \sum_{k=p+1}^n \theta_k^2 \right] + \frac{1}{\beta^2} \frac{1}{\sigma_{min}^2}} \\ &= ||A^{\dagger}|| \sqrt{\frac{1}{\alpha^2} \left( ||x||^2 + ||A^{\dagger}||^2 ||r||^2 \right) + \frac{1}{\beta^2}}, \end{split}$$

where  $A^{\dagger}$  denotes the Moore–Penrose inverse (see [9, p. 421]) of A, and x denotes the solution of the linear least squares problem associated with A and b.

Looking at the general result of Proposition 3.3, we see that the quantities (scalars) involved in the computation of the  $x_r$  condition number are nothing but the singular values  $\sigma_k$  of A and the components  $\theta_k$  of b along singular vectors  $u_k$ . Finally, observe that the critical gap is  $\sigma_r - \sigma_{r+1}$ .

4. Numerical experiments. We now describe some numerical tests carried out in MATLAB. Our test cases come from the package Regularization Tools<sup>4</sup> by Hansen [7]. We arbitrarily choose values of n, p, and r. By means of a specific routine in this package, we generate pairs (A, b) associated with some test problems indicated by their name. To validate the expression of the exact condition number, we use the numerical derivative code<sup>5</sup> by D'Errico, called jacobianest.m. to estimate the corresponding Jacobian at a given particular point z. The code Jacobianest.m

<sup>&</sup>lt;sup>4</sup>See http://www2.imm.dtu.dk/~pch/Regutools.

<sup>&</sup>lt;sup>5</sup>See http://www.mathworks.com/matlabcentral/fileexchang/13490-automatic-numericaldifferentiation.

the estimates to sixth order. For our purpose, we recast  $x_r$  as  $f: z = vec([A, b]) \to x_r$ prior to the use of jacobianest.m, and then we approximate the condition number of  $x_r$  as the 2-norm of the estimated Jacobian. Note that in all tests we set  $\alpha = \beta = 1$ .

Table 1 displays the exact condition number versus an estimate of the condition number produced with jacobianest.m. The results show how the derived expression of the exact condition compares well with the finite difference estimate.

5. Conclusion. We solved the problem of the determination of a closed formula for the condition number of the truncated singular value solution of an ill-posed problem which relies on a singular value decomposition of the problem. We anticipate that the presented formula will therefore stimulate research in several directions. Finding good estimates of the condition number using iterative techniques would, for instance, be of crucial relevance for large-scale problems. From a theoretical point of view, we also believe that the condition number may bring new insight into the problem of the detection of the truncation index of the singular value decomposition. One of the topics of future research will be to explore this issue for practical problems.

Acknowledgment. The authors thank the anonymous referees for the valuable comments and suggestions.

### REFERENCES

- [1] Å. BJÖRCK, Numerical Methods for Least Squares Problems, SIAM, Philadelphia, 1996.
- [2] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, 3rd ed., Johns Hopkins University Press, Baltimore, MD, 1996.
- [3] S. GRATTON, On the condition number of linear least squares problems in Frobenius norm, BIT, 36 (1995), pp. 523–530.
- [4] S. GRATTON, D. TITLEY-PELOQUIN, AND J. TSHIMANGA, Sensitivity and conditioning of the truncated total least squares solution, SIAM J. Matrix. Anal. Appl., 34 (2013), pp. 1257–1276.
- [5] P. C. HANSEN, The truncated SVD as a method for regularization, BIT, 27 (1987), pp. 534–553.
- P. C. HANSEN, Rank-Deficient and Discrete Ill-Posed Problems: Numerical Aspects of Linear Inversion, SIAM, Philadelphia, 1998.
- [7] P. C. HANSEN, Regularization tools: A MATLAB package for analysis and solution of discrete ill-posed problems, Numer. Algorithms, 46 (2007), pp. 187–194.
- [8] N. J. HIGHAM, Accuracy and Stability of Numerical Algorithms, 2nd ed., SIAM, Philadelphia, 2002.
- [9] R. A. HORN AND C. R. JOHNSON, *Matrix Analysis*, Cambridge University Press, Cambridge, UK, 1990.
- [10] R. A. HORN AND C. R. JOHNSON, Topics in Matrix Analysis, Cambridge University Press, Cambridge, UK, 1991.
- [11] Z.-C. LI, H.-T. HUANG, AND Y. WEI, Ill-conditioning of the truncated singular value decomposition, Tikhonov regularization and their applications to numerical partial differential equations, Numer. Linear Algebra Appl., 18 (2011), pp. 205–221.
- [12] G. W. STEWART, Error and perturbation bounds for subspaces associated with certain eigenvalue problems, SIAM Rev., 15 (1973), pp. 727–764.
- [13] G. W. STEWART AND J.-G. SUN, *Matrix Perturbation Theory*, Academic Press, New York, 1990.
- [14] C. R. VOGEL, Computational Methods for Inverse Problems, Frontiers in Appl. Math., SIAM, Philadelphia, 2002.