

An orthogonal high relative accuracy algorithm for the symmetric eigenproblem*

Froilán M. Dopico, Juan M. Molera and Julio Moro[†]
Departamento de Matemáticas
Universidad Carlos III de Madrid
28911-Leganés, Spain

May 7, 2003

Abstract

We propose a new algorithm for the symmetric eigenproblem which computes eigenvalues and eigenvectors with high relative accuracy for the largest class of symmetric, definite and indefinite, matrices known so far. The algorithm is divided in two stages: the first one computes a SVD with high relative accuracy, and the second one obtains eigenvalues and eigenvectors from singular values and vectors. Using the SVD as a first stage is responsible both for the wide applicability of the algorithm and for being able to use only orthogonal transformations, unlike previous algorithms in the literature. Theory, a complete error analysis and numerical experiments are presented.

Key identifying words: symmetric eigenproblem, singular value decomposition, high relative accuracy.

AMS subject classification. 65F15, 65G50, 15A18.

1 Introduction

An *orthogonal spectral decomposition* of a real symmetric n by n matrix A is a factorization $A = Q \Lambda Q^T$, where Q is real orthogonal and $\Lambda = \text{diag}[\lambda_1, \dots, \lambda_n]$ is diagonal. We assume that $\lambda_1 \geq \dots \geq \lambda_n$. The columns q_i , $i = 1, \dots, n$, of Q are the eigenvectors of A corresponding to the eigenvalues λ_i , $i = 1, \dots, n$. In this paper we present an algorithm that computes an orthogonal spectral decomposition for the largest class so far of symmetric matrices with the following *high relative accuracy*:

*The research conducting to this paper was partially supported by the spanish Ministerio de Ciencia y Tecnología through project BFM-2000-0008.

[†]emails: dopico@math.uc3m.es, molera@math.uc3m.es, jmorero@math.uc3m.es

- The error in each computed eigenvalue, $\widehat{\lambda}_i$, is

$$|\lambda_i - \widehat{\lambda}_i| = O(\kappa \epsilon) |\lambda_i|, \quad (1)$$

where we assume that $\widehat{\lambda}_1 \geq \dots \geq \widehat{\lambda}_n$, ϵ is the unit roundoff of the finite arithmetic employed in the computation and κ is a relevant condition number, to be specified later in §2.1. High relative accuracy will be achieved whenever $\kappa \approx 1$.

- The angle $\Theta(q_i, \widehat{q}_i)$ between each computed eigenvector \widehat{q}_i and the exact one q_i satisfies

$$\Theta(q_i, \widehat{q}_i) = \frac{O(\kappa \epsilon)}{\text{relgap}(|\lambda_i|)}, \quad (2)$$

where $\text{relgap}(|\lambda_i|) = \min \left\{ \min_{j \neq i} \left| \frac{|\lambda_j| - |\lambda_i|}{|\lambda_i|} \right|, 1 \right\}$. This may be a considerably more accurate outcome than that of a conventional eigenvalue method, like QR, divide-and-conquer or bisection with inverse iteration, for instance. Such algorithms produce results with high *absolute* accuracy, i.e. satisfying

$$|\lambda_i - \widehat{\lambda}_i| = O(\epsilon) \max_j |\lambda_j|,$$

instead of (1), and

$$\Theta(q_i, \widehat{q}_i) = \frac{O(\epsilon)}{\frac{\min_{j \neq i} |\lambda_i - \lambda_j|}{\max_j |\lambda_j|}},$$

instead of (2). Thus, a conventional algorithm may provide approximations for the small eigenvalues ($\frac{\max_j |\lambda_j|}{|\lambda_i|} \sim \frac{1}{\epsilon}$) with no correct significant digits, or even with the wrong sign. Moreover, if there are two or more small eigenvalues, their eigenvectors may be computed very inaccurately, even when the eigenvalues are well-separated in the relative sense (e.g. $\lambda_i = 10^{-15}$ and $\lambda_j = 10^{-16}$ if $\lambda_1 = 1$). At present, high relative accuracy can only be reached for certain classes of *symmetric* matrices.

Identifying classes of matrices for which either a *singular value decomposition* (SVD) or a spectral decomposition can be computed with high relative accuracy has been a very active area of research in the last fifteen years (see [7] and references therein for an overview). At present, high relative accuracy eigensolvers are only available for some symmetric matrices, and are far less developed than accurate SVD algorithms (except, of course, in the related positive definite case [8]). To be more precise, several easily characterized classes of matrices have been identified in [7] for which high relative accuracy SVDs can be computed, while present symmetric indefinite eigensolvers deliver high relative accuracy for matrices which are not easy to recognize (with the exception of scaled diagonally dominant matrices [2]). As can be seen in [23, 28], the symmetric indefinite matrices deserving high relative accuracy spectral decompositions have been characterized through the positive semidefinite polar factor, which is difficult to compute. In this regard, *the main contribution of the present paper is to prove that the proposed eigensolver achieves high relative accuracy (1),(2) for all symmetric matrices in any of the classes identified in [7]*. Moreover, it will do so, under very general assumptions, for any class of matrices eventually identified in the future for which high relative accuracy SVDs can be computed. None of the present symmetric eigensolvers can guarantee, to our knowledge, high relative accuracy for the classes of matrices above.

The most general algorithm to compute spectral decompositions of symmetric indefinite matrices with high relative accuracy is at present the so-called *implicit J-orthogonal* algorithm. It was introduced by Veselić in [27] and carefully analyzed by Slapničar in [23]. It has the

advantage, when delivering high relative accuracy, of producing an error smaller than (2) in the eigenvectors, namely

$$\Theta(q_i, \hat{q}_i) = \frac{O(\tilde{\kappa}\epsilon)}{\text{relgap}(\lambda_i)} \quad (3)$$

for an associated condition number $\tilde{\kappa}$, which has been observed in practice to be $\tilde{\kappa} \approx 1$. This error depends on

$$\text{relgap}(\lambda_i) = \min \left\{ \min_{j \neq i} \frac{|\lambda_j - \lambda_i|}{|\lambda_i|}, 1 \right\}, \quad (4)$$

the relative gap between the eigenvalues, the natural one for the symmetric eigenvalue problem, usually larger than the relative gap between *the absolute values* of the eigenvalues.

There are other important differences between the algorithm by Veselić and Slapničar and the one proposed below: the J-orthogonal algorithm uses *hyperbolic* transformations [16, §12.5.4], which complicates the error analysis and increases the constants in the error bounds. The algorithm we propose here uses only *orthogonal* transformations. Also, the error bounds for the hyperbolic J-orthogonal algorithm are valid modulo a minor proviso (bounded growth of the scaled condition number of certain matrices appearing in each step of the iteration), while the new algorithm can be implemented in such a way that no proviso is needed to guarantee the error bounds. On the other hand, the J-orthogonal algorithm may be easily extended to matrix pencils, while this is not possible for the one presented here. There are also similarities: both algorithms require to previously factorize the matrix, and both depend crucially on employing algorithms of one-sided Jacobi type.

The basic motivation for the algorithm we propose is to take advantage of the intense research effort, undertaken in the last few years, focused on identifying classes of matrices for which an SVD can be computed with high relative accuracy. This collective effort has led to the unified approach in [7]. The connection with our work lies in that the SVD and the spectral decomposition are closely related for symmetric¹ matrices: the singular values are the absolute values of the eigenvalues, and eigenvectors may be constructed from singular vectors. To be more precise, let $A = U\Sigma V^T$ be an SVD of $A = A^T$, where U, V are n by n orthogonal with columns u_i, v_i , $i = 1, \dots, n$, and $\Sigma = \text{diag}[\sigma_1, \dots, \sigma_n]$ with $\sigma_1 \geq \dots \geq \sigma_n \geq 0$. In the simplest (and most frequent) case in which all singular values of A are distinct, the eigenvalues of A are

$$(v_i^T u_i) \sigma_i, \quad (5)$$

with $v_i^T u_i = \pm 1$ for all $i = 1, \dots, n$, and the corresponding eigenvectors are v_i (u_i may be equally chosen). Hence, once an SVD is known, the only additional work to obtain the eigenvalues is to determine the sign ± 1 via the scalar product $v_i^T u_i$ of right and left singular vectors (the general case when groups of equal singular values appear is presented in §3.1). Notice that $v_i^T A v_i = v_i^T u_i \sigma_i$, i.e. the scalar product above can be thought of as a cheaper and indirect way of obtaining the sign from the Rayleigh quotient, avoiding the multiplication by the matrix A , which may give the wrong sign due to its large condition number (one example of this difficulty will be shown at the end of §3.3). In fact, this particular way of assigning the signs through $v_i^T u_i$, together with the proof of its accuracy, is one of the crucial issues in this paper.

Therefore, given a computed high relative accuracy SVD of $A = A^T$ satisfying

¹All the results in this paper are valid for Hermitian matrices, although for the sake of simplicity we restrict the discussion to the real symmetric case.

$$|\sigma_i - \hat{\sigma}_i| = O(\kappa\epsilon) |\sigma_i| \quad (6)$$

$$\Theta(v_i, \hat{v}_i) = \frac{O(\kappa\epsilon)}{\text{relgap}(\sigma_i)}, \quad \Theta(u_i, \hat{u}_i) = \frac{O(\kappa\epsilon)}{\text{relgap}(\sigma_i)}, \quad (7)$$

with

$$\text{relgap}(\sigma_i) = \min \left\{ \min_{j \neq i} \frac{|\sigma_i - \sigma_j|}{\sigma_i}, 1 \right\}, \quad (8)$$

it only remains to prove that the pair \hat{v}_i, \hat{u}_i approximates the pair v_i, u_i closely enough so that the computed value of the scalar product approximates ± 1 with an *absolute* error smaller than one, to get (1) and (2). Notice that this is no longer a high relative accuracy problem.

In this spirit we propose the following algorithm to compute the eigenvalues and eigenvectors of a symmetric matrix in two stages:

Stage 1: Compute an SVD of A with accuracy (6) and (7).

Stage 2: Compute the eigenvalues of A by giving signs, according to (5), to the singular values computed in **Stage 1**. The corresponding eigenvectors are the right (or left) singular vectors computed in **Stage 1**. When groups of equal singular values are present, this step becomes more involved (see §3.3 below).

We will show that Stage 2 provides high relative accuracy in the eigenvalues (1) and in the eigenvectors (2) as long as Stage 1 gives an SVD with small backward multiplicative error (as in formula (16) below, that in turn guarantees (6) and (7)). As to Stage 1, there are at present algorithms to perform it for several classes of matrices, summarized in [7]. These are the algorithms we are going to use, although any future high relative accuracy SVD algorithm may be employed for Stage 1.

Maybe the most remarkable contribution of Demmel et al. in [7] is developing algorithms which compute high relative accuracy SVDs (i.e., satisfying (6) and (7)) for *any* matrix such that a so-called *rank-revealing decomposition* (RRD) can be computed with enough accuracy. A RRD of $G \in \mathbb{R}^{m \times n}$, $m \geq n$, is a factorization $G = \mathcal{X}\mathcal{D}\mathcal{Y}^T$ with $\mathcal{D} \in \mathbb{R}^{r \times r}$ diagonal and nonsingular and $\mathcal{X} \in \mathbb{R}^{m \times r}$, $\mathcal{Y} \in \mathbb{R}^{n \times r}$, where both matrices \mathcal{X}, \mathcal{Y} have full column rank and are well-conditioned (notice that this implies $r = \text{rank}(G)$). According to the structure of the algorithms in [7], a more specific description of the *SSVD algorithm* we propose here (the *SSVD* stands for *Signed SVD*) is the following

Algorithm 1 (SSVD)

Input: Symmetric matrix A .

Output: Eigenvalues $\Lambda = \text{diag}[\lambda_i]$ and eigenvectors $Q = [q_1 \dots q_n]$; $A = Q\Lambda Q^T$.

1. Compute a RRD factorization XDY^T of A .
2. Compute SVD $XDY^T = USV^T$ of RRD using algorithms in [7, §3].
3. Compute the eigenvalues of A , by giving signs to the singular values, and the corresponding eigenvectors, using either Algorithm 2 (SYSSV) in page 12, or Algorithm 3 (SYSSVR) in page 25 below.

The accuracy required in [7] on the *computed* RRD matrices X, D, Y in order to guarantee that a high relative accuracy SVD can be obtained is given by the following forward error bounds:

$$\begin{aligned} |d_{ii} - \delta_{ii}| &= O(\epsilon)|\delta_{ii}|, \\ \|X - \mathcal{X}\| &= O(\epsilon)\|\mathcal{X}\|, \\ \|Y - \mathcal{Y}\| &= O(\epsilon)\|\mathcal{Y}\|, \end{aligned} \tag{9}$$

where $\|\cdot\|$ denotes the spectral norm and d_{ii}, δ_{ii} stand, respectively, for the diagonal elements of D, \mathcal{D} . Once a RRD factorization XDY^T satisfying (9) is available, either Algorithm 3.1 or Algorithm 3.2 of [7] provide a high relative accuracy SVD of XDY^T with overall relative error (including the initial factorization stage) of order $O(\epsilon \max\{\kappa(X), \kappa(Y)\})$ in the singular values, and $O(\epsilon \max\{\kappa(X), \kappa(Y)\})$ over the relative gap (8) in the singular vectors, where $\kappa(\cdot)$ denotes the condition number in the spectral norm. The key to prove this high relative accuracy is that both the error (9) in the factorization and the errors introduced by Algorithms 3.1 or 3.2 of [7] produce a backward error of multiplicative type, instead of the additive type usually produced by conventional algorithms (see §2.1 for a more detailed discussion).

Several classes of matrices have been found in the last ten years for which it is possible to compute an accurate RRD. They include bidiagonal, acyclic, Cauchy, Vandermonde, graded, scaled diagonally dominant matrices, as well as all well-scalable symmetric positive definite and some well-scalable symmetric indefinite matrices, etc. Hence, for all symmetric matrices in any of the classes described in [7, p. 26-27], Algorithm 1 will produce a spectral decomposition with the high relative accuracy given by (1) and (2) under the criteria given in [7] for computing accurate RRDs.

The way RRDs with the accuracy (9) are obtained in [7] is the reason why the J-orthogonal algorithm does not guarantee, at present, high relative accuracy for the classes of symmetric matrices mentioned above. The key is that RRDs are computed in [7] from Gaussian elimination with complete pivoting (GECP)². Moreover, a plain implementation of GECP does not guarantee accuracy (9) for most of the classes in [7]. This can only be achieved through special, nontrivial implementations of GECP, sometimes demanding a great deal of ingenuity (see [7, 6]). Since GECP leads in general to RRDs with $X \neq Y$, even if the matrix to be factorized is symmetric, the J-orthogonal algorithm cannot be directly applied because it begins with a more restrictive *symmetric indefinite factorization* SJS^T of the matrix $A = A^T$, where J is square diagonal with diagonal elements ± 1 , and S has full column rank³. At present, for instance, it is not known whether some modifications in the algorithm leading to the symmetric indefinite factorization would ensure that it is accurately computed in the sense of (9) for the classes of symmetric matrices in [7]. Although [24] contains a thorough error analysis of the symmetric indefinite factorization, no general forward error bounds are provided allowing to determine classes of matrices admitting an accurate factorization.

Notice that the nonsymmetric character of Algorithm **SSVD** is responsible both for making it valid for a large class of matrices and for being able to use only orthogonal transformations in **step 2**, while the J-orthogonal algorithm has to use hyperbolic ones. The prize to pay, however, is that by applying an SVD algorithm (valid for any matrix) to a symmetric matrix, we are not making any use of the symmetry of A . Thus, the algorithm is not backward stable, in the sense that one cannot guarantee that the computed eigenvalues and eigenvectors are the exact

²Some mention is also made in [7] of using QR with complete pivoting [22, 18, 4]. This would open the possibility of using Algorithm 3.3 of [7], less costly than Algorithms 3.1-3.2 for **step 2** of **SSVD**.

³Notice that, although SJS^T is not an RRD, its computation is equivalent to computing a symmetric RRD of the form XDY^T , see [24].

eigenvalues and eigenvectors of a close *symmetric* matrix. An important consequence of this is that a first, plain implementation of the algorithm produces an error (2) in the eigenvectors depending on the relative gap between the absolute values of the eigenvalues, i.e. between the singular values. This does not appear if we use a symmetric algorithm (as the J-orthogonal algorithm) producing a *symmetric* backward error, since in that case the relative perturbation theory for symmetric matrices [15, 20, 28] leads to eigenvector error bounds depending on the natural eigenvalue relative gap (3).

This is the reason why we analyze two different implementations of **step 3** of Algorithm 1. The first and simplest one, Algorithm 2, follows straightforwardly the ideas explained after equation (5) and delivers the announced accuracy (1), (2). To improve the accuracy (2) for the eigenvectors we offer in section 5 a second, more sophisticated implementation, Algorithm 3. Although this is the one we recommend in practice, we also include (and analyze) the first implementation because it is very difficult to understand how and why Algorithm 3 works without the error analysis for Algorithm 2. We stress that both versions compute the same eigenvalues, and only differ in the eigenvector computation step, which is more accurate for Algorithm 3. To be more precise, notice that the relative gap between the singular values, $relgap(\sigma_i)$, can be much smaller than the relative gap between the eigenvalues only when $relgap(\sigma_i)$ is calculated from two consecutive singular values coming from eigenvalues of different sign. We show in Section 5 that if λ_k is the eigenvalue corresponding to σ_i , and λ_j is the eigenvalue corresponding to the singular value closest to σ_i , then the errors $\Theta(q_j, \hat{q}_j)$, $\Theta(q_k, \hat{q}_k)$ in both eigenvectors computed by Algorithm 1, using Algorithm 3 in **step 3**, are bounded by

$$\frac{O(\kappa \epsilon)}{\min\{relgap(\lambda_k), relgap(\lambda_j)\}} \quad (10)$$

instead of (2), i.e., *given two close singular values coming from eigenvalues of different sign, the error in the corresponding eigenvectors does not depend on the distance between these singular values*. Actually, we will prove that the accuracy (10) is already achieved by Algorithm 2 in all cases except in those with an eigenvalue sign distribution like the one mentioned above. And even in those cases, Algorithm 2 achieves (10) if $relgap(\sigma_i) \leq \kappa \epsilon$. Notice that the worst eigenvector error bound for Algorithm 1 using Algorithm 3 for **step 3** is also reached by the J-orthogonal algorithm, although some of the eigenvectors computed by the J-orthogonal algorithm can be more accurate than the corresponding ones computed by SSVD.

Concerning the computational cost of Algorithm 1, it is $O(n^3)$ provided the initial RRD costs $O(n^3)$ (some classes of matrices allow an accurate RRD, but not at $O(n^3)$ cost [7]). As usual for high accuracy algorithms, Algorithm 1 is more expensive than other $O(n^3)$ conventional eigenvalue methods, like QR, divide-and-conquer, etc. The most expensive part of Algorithm SSVD is the one-sided Jacobi method employed in **step 2**. However, some ways have been recently found [13] to speed up one-sided Jacobi, about as fast as the QR algorithm for SVD.

It is difficult to compare the cost of Algorithm 1 with that of the J-orthogonal algorithm. If in both cases we do not count the initial factorization, the difference between Algorithm 3.1 of [7] and Algorithm 3.3.1 of [23] seems to amount to two matrix multiplications and one QR factorization. However, numerical experience indicates that Algorithm 3.1 of [7] requires less Jacobi sweeps than Algorithm 3.3.1 of [23] (see section 6.2). Finally, **step 3** of Algorithm SSVD costs, in general, $O(n^2)$, but for every cluster with d close singular values corresponding to eigenvalues of both signs, and if eigenvectors need to be computed, there is an overhead cost of $O(d^3) + O(d^2n)$. Clearly, this is maximized when only one cluster of size $d = n$ is present. Then, the cost of **step 3** is $O(n^3)$. As to the timing statistics, the run-times of both algorithms will be shown to be comparable in the numerical experiments below.

The rest of the paper is organized as follows: Section 2 collects the mathematical results required to perform a complete error analysis of Algorithm 1. Section 3 describes in detail Algorithm 2 for **step 3** of Algorithm 1 (SSVD), including the corresponding pseudocode. Special attention is given to the presence of clusters of close singular values. Section 4 contains a detailed error analysis of this implementation of Algorithm 1 in the most general setting, allowing for the presence of clusters. Actually, the possible presence of clusters is the reason why we need to devote a complete section to the error analysis, which is straightforward in the case of matrices with distinct singular values. The error analysis in Section 4 motivates Algorithm 3 for **step 3** of Algorithm 1, which is developed and analyzed in section 5, where (10) is proved in the most general setting, with any distribution of clusters. Section 6 addresses the practical implementation of Algorithm 1, together with the numerical tests. Conclusions and discussion of open problems are presented in Section 7. Finally, two appendices contain certain proofs of results in, respectively, sections 2 and 5 which, due to their length and technical character, are better postponed to the end of the paper.

2 Preliminary results

We collect in this section the mathematical results required to perform the error analysis of Algorithm 1. As stated in the introduction, the only requirement on the high relative accuracy SVD algorithm in **step 2** of Algorithm 1 is producing a small multiplicative backward error when performed in finite arithmetic. A precise statement is given in §2.1 for algorithms in [7]. It is also shown that the error due to the initial RRD can be absorbed as an additional multiplicative backward error. Section 2.2 summarizes the multiplicative perturbation theory for singular values and for bases of singular subspaces needed to guarantee the high relative accuracy of the overall algorithm.

2.1 Backward error of the SVD algorithm

The following theorem is essentially proved in [7]

Theorem 2.1 *Algorithm 3.1 of [7] (see Algorithm 4 in § 6.1 below) produces a multiplicative backward error when executed with machine precision ϵ , i.e. if $G = XDY^T \in \mathbb{R}^{m \times n}$ is the RRD computed in **step 1** of Algorithm 1 and $\widehat{U}\widehat{\Sigma}\widehat{V}^T$ is the SVD computed by the algorithm, then there exist matrices $U' \in \mathbb{R}^{m \times r}$, $V' \in \mathbb{R}^{n \times r}$, $E \in \mathbb{R}^{m \times m}$, $F \in \mathbb{R}^{n \times n}$ such that U' and V' have orthonormal columns,*

$$\begin{aligned} \|U' - \widehat{U}\| &= O(\epsilon), & \|V' - \widehat{V}\| &= O(\epsilon), \\ \|E\| &= O(\epsilon\kappa(X)), & \|F\| &= O(\epsilon\kappa(R')\kappa(Y)), \end{aligned} \tag{11}$$

where R' is the best conditioned row diagonal scaling of the triangular matrix R appearing in **step 1** of Algorithm 3.1 of [7] (i.e. Algorithm 4 in § 6.1) and

$$(I + E)G(I + F) = U'\widehat{\Sigma}V'^T \tag{12}$$

It is proved in [7] that $\kappa(R')$ is at most of order $O(n^{3/2}\kappa(X))$, but in practice we have observed in extensive numerical tests that $\kappa(R')$ behaves as $O(n)$. One can get rid of the factor $\kappa(R')$ at the prize of using the more costly Algorithm 3.2 of [7].

We state Theorem 2.1 because the original result, [7, Thm. 3.1], is not phrased as a backward error result, which is what we need for the subsequent error analysis. The only missing piece in

the analysis of [7] is the fact that one-sided Jacobi [16, §8.6.3] produces a small multiplicative backward error. This can be easily derived from Proposition 3.13 in [12] and, since it is not central to our argument, we defer its proof, together with that of Theorem 2.1, to Appendix A below. Two different versions of Algorithm 4 will be analyzed in the appendix, depending on whether the right-handed or left-handed version of one-sided Jacobi is employed. We will show that the right-handed version, i.e. the one in which the Jacobi rotations are applied from the right, guarantees smaller error bounds and leads precisely to Theorem 2.1. For the left-handed version one can prove a theorem similar to 2.1, but with a weaker error bound for F , and requiring a minor proviso to ensure the accuracy. However, the left-handed version is still the one usually employed in practice since it is much faster and no significant difference has ever been observed in accuracy. This is why we use it in most of the experiments in section 6. Finally, it is crucial for the accuracy of one-sided Jacobi algorithms to impose as stopping criterion that the cosines of the angles between the different columns (or rows, depending on the version of one-sided Jacobi) are smaller than ϵ times the dimension of the matrix.

Once the backward error of the SVD algorithm is shown to be multiplicative, the perturbation theory in §2.2 below can be used to prove high relative accuracy, namely that the computed singular values and vectors of XDY^T satisfy

$$\begin{aligned} |\sigma_i - \hat{\sigma}_i| &= O(\kappa \epsilon) \sigma_i, \\ \Theta(v_i, \hat{v}_i) &= \frac{O(\kappa \epsilon)}{\text{relgap}(\sigma_i)}, \\ \Theta(u_i, \hat{u}_i) &= \frac{O(\kappa \epsilon)}{\text{relgap}(\sigma_i)}. \end{aligned} \tag{13}$$

where

$$\kappa = \kappa(R') \max\{\kappa(X), \kappa(Y)\} \tag{14}$$

is in this case the relevant condition number announced in the introduction.

As a matter of fact, one may even absorb into a backward error of the form (12) the error in the initial RRD, i.e. the one due to the fact that the SVD computation does not start from the symmetric matrix A itself, but from its computed rank-revealing decomposition: let $A = \mathcal{X}\mathcal{D}\mathcal{Y}^T$ be an exact RRD factorization of A and assume the starting decomposition XDY^T has been computed accurately enough so that the computed matrices X, D, Y satisfy conditions (9). Then, as shown in the proof of Theorem 2.1 in [7], there exist matrices E_f, F_f with $\|E_f\| = O(\epsilon\kappa(X))$, $\|F_f\| = O(\epsilon\kappa(Y))$ such that

$$(I + E_f)A(I + F_f) = XDY^T. \tag{15}$$

This, together with (12), implies that

$$U'\hat{\Sigma}V'^T = (I + \tilde{E})A(I + \tilde{F}) \tag{16}$$

where the backward errors \tilde{E}, \tilde{F} are of size $\|\tilde{E}\| = O(\epsilon\kappa(X))$, $\|\tilde{F}\| = O(\epsilon\kappa(R')\kappa(Y))$ and reflect that the errors produced by both the RRD factorization and the SVD algorithm are backward multiplicative.

We stress that all our error analysis is done in terms of the backward errors $\|\tilde{E}\|$ and $\|\tilde{F}\|$. Although we have focused on the case when $\|E_f\| = O(\epsilon\kappa(X))$ and $\|F_f\| = O(\epsilon\kappa(Y))$, any other more general backward errors for the factorization step can be trivially incorporated into the error analysis, since, up to first order,

$$\|\tilde{E}\| \leq \|E_f\| + O(\epsilon\kappa(X)), \quad \|\tilde{F}\| \leq \|F_f\| + O(\epsilon\kappa(R')\kappa(Y)).$$

2.2 Multiplicative perturbation theory

Let G be a real m by n matrix with SVD $G = U\Sigma V^T$ and singular values $\sigma_1 \geq \sigma_2 \geq \dots$. We consider a multiplicative perturbation $\tilde{G} = (I + E)G(I + F)$ of G with SVD $\tilde{G} = \tilde{U}\tilde{\Sigma}\tilde{V}^T$ and singular values $\tilde{\sigma}_i$, also ordered decreasingly.

Theorem 2.2 [15, Theorem 3.1] *Let $G \in \mathbb{R}^{m \times n}$, $\tilde{G} = (I + E)G(I + F)$, and set*

$$\eta = \max\{\|E\|, \|F\|\}, \quad \eta' = 2\eta + \eta^2. \quad (17)$$

Then

$$\frac{|\sigma_i - \tilde{\sigma}_i|}{\sigma_i} \leq \eta'.$$

Besides the change in the singular values, we also need to estimate the changes undergone by singular subspaces or, more precisely, by their bases. Although the following results are valid for rectangular matrices (see [20, 9]), we state them in the square case, the only one we deal with in Section 4. Thus, G is now a real n by n matrix and $\tilde{G} = (I + E)G(I + F)$. Let

$$G = [U_1 \ U_2] \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix} \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix}, \quad (18)$$

$$\tilde{G} = [\tilde{U}_1 \ \tilde{U}_2] \begin{bmatrix} \tilde{\Sigma}_1 & 0 \\ 0 & \tilde{\Sigma}_2 \end{bmatrix} \begin{bmatrix} \tilde{V}_1^T \\ \tilde{V}_2^T \end{bmatrix} \quad (19)$$

be two conformally partitioned SVDs of G and \tilde{G} , i.e. the four matrices $\Sigma_1, \tilde{\Sigma}_1 \in \mathbb{R}^{q \times q}$ and $\Sigma_2, \tilde{\Sigma}_2 \in \mathbb{R}^{(n-q) \times (n-q)}$ are diagonal. No particular order is assumed on the singular values. The change in the singular subspaces is usually measured through the sines of the canonical angles $\Theta(U_1, \tilde{U}_1)$ between the column spaces of U_1 and \tilde{U}_1 , and $\Theta(V_1, \tilde{V}_1)$ between the column spaces of V_1 and \tilde{V}_1 (see [26]). It is well-known that this change is governed (see e.g. [20, Thm. 4.1]) by the singular value relative gap

$$rg(\Sigma_1, \tilde{\Sigma}_2) = \min_{\substack{\sigma \in \sigma(\Sigma_1) \\ \tilde{\sigma} \in \sigma(\tilde{\Sigma}_2)}} \frac{|\sigma - \tilde{\sigma}|}{\tilde{\sigma}}, \quad (20)$$

where $\sigma(M)$ denotes the set of singular values of the matrix M .

This kind of results, however, is not enough for our purposes. The fact that the signs of the eigenvalues are obtained through scalar products like the one in (5) forces us to accurately compute not only the singular subspaces, but also the corresponding *simultaneous* bases U_i and V_i . To ensure this, finer perturbation results are needed, dealing with the sensitivity of the bases themselves. It has been observed in [9] that simultaneous bases of singular subspaces do not have the same sensitivity under perturbation as their corresponding singular subspaces. More precisely, bases may be much more sensitive to *additive* perturbations than singular subspaces. Fortunately enough for our purposes, both sensitivities are essentially equal for multiplicative perturbations. A detailed discussion of these issues may be found in [9, 10], including a stronger version of the following result (we use the Frobenius norm $\|\cdot\|_F$, as usual when the dimension of the subspaces is larger than 1).

Theorem 2.3 [9, Theorem 2.2] *Let $G \in \mathbb{R}^{n \times n}$ and $\tilde{G} = (I + E)G(I + F)$ with respective SVDs (18) and (19). Then there exists an orthogonal matrix $P \in \mathbb{R}^{q \times q}$ such that*

$$\sqrt{\|U_1 P - \tilde{U}_1\|_F^2 + \|V_1 P - \tilde{V}_1\|_F^2} \leq 2\sqrt{q} \left[\eta + \frac{\eta'}{1 - \eta} \frac{1}{\text{relgap}(\Sigma_1, \tilde{\Sigma}_2)} \right], \quad (21)$$

where $\text{relgap}(\Sigma_1, \tilde{\Sigma}_2)$ is given by

$$\text{relgap}(\Sigma_1, \tilde{\Sigma}_2) = \min\{rg(\Sigma_1, \tilde{\Sigma}_2), 1\}. \quad (22)$$

and η, η' are given by (17).

Although it is more usual in the literature [7, 6] to define the relative gap (20) with the roles of Σ_1 and Σ_2 reversed, we have chosen the definition above to conform to the cited perturbation theorems. However, this does not represent any significant difference in the error bounds, since a straightforward calculation shows that

$$2 \text{relgap}(\tilde{\Sigma}_2, \Sigma_1) \geq \text{relgap}(\Sigma_1, \tilde{\Sigma}_2) \geq \frac{1}{2} \text{relgap}(\tilde{\Sigma}_2, \Sigma_1). \quad (23)$$

On the other hand, as usual in this kind of perturbation bounds, one can reformulate the definition of the gaps in order to make them depend only on the unperturbed singular values, at the cost of somewhat complicating the bounds.

The main point of Theorem 2.3 is that the orthogonal matrix P is the same for both left and right singular vectors. This will be enough to guarantee the accuracy of the sign assignment and of the computed bases of invariant subspaces⁴.

3 Computing spectral decompositions from SVDs

This section is divided in three parts: §3.1 outlines the mathematical basis for the main idea underlying Algorithm 1, namely that one can easily get a spectral decomposition of a symmetric matrix if one is given an SVD, even if the matrix has groups of equal singular values. Some practical details concerning clusters of close singular values in finite precision will be considered in §3.2. The complete pseudocode for Algorithm 2, the simplest implementation of **step 3** in Algorithm 1, will be presented in §3.3.

3.1 Mathematical basis

Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix with SVD $A = U\Sigma V^T$. Then, $V^T A V = V^T U \Sigma$ is orthogonally similar to A with $V^T U$ orthogonal. If A has distinct singular values $\sigma_1 > \sigma_2 > \dots > \sigma_p$ with respective multiplicities $m_i, i = 1, \dots, p$ ($m_1 + \dots + m_p = n$) and we partition U and V accordingly as

$$U = [\mathcal{U}_1 \mid \mathcal{U}_2 \mid \dots \mid \mathcal{U}_p],$$

$$V = [\mathcal{V}_1 \mid \mathcal{V}_2 \mid \dots \mid \mathcal{V}_p]$$

into blocks $\mathcal{U}_i, \mathcal{V}_i \in \mathbb{R}^{n \times m_i}$ corresponding to each distinct singular value, then

$$\mathcal{V}_i^T \mathcal{U}_j = 0 \quad \text{whenever } i \neq j \quad (24)$$

⁴Actually, Theorem 2.3 is stronger than the usual bounds on the canonical angles between singular subspaces, since one can easily show that $\|\sin(\Theta(U_1, \tilde{U}_1))\|_F \leq \|U_1 P - \tilde{U}_1\|_F$, and similarly for V_1 .

since, due to the symmetry of A , both its left and its right singular vectors are eigenvectors of A^2 . Consequently,

$$V^T U = \text{diag}[\mathcal{V}_1^T \mathcal{U}_1, \dots, \mathcal{V}_p^T \mathcal{U}_p], \quad (25)$$

is block-diagonal, where each diagonal block $\mathcal{V}_i^T \mathcal{U}_i \in \mathbb{R}^{m_i \times m_i}$ is itself orthogonal. Furthermore, since

$$V^T A V = \text{diag}[\sigma_1 \mathcal{V}_1^T \mathcal{U}_1, \dots, \sigma_p \mathcal{V}_p^T \mathcal{U}_p] \quad (26)$$

is symmetric, we conclude that each $\mathcal{V}_i^T \mathcal{U}_i$ is not only orthogonal, but also symmetric and its eigenvalues, ± 1 , are precisely the signs of those eigenvalues of A having modulus σ_i . In the simplest case when $m_i = 1$, the eigenvalue is just $v_i^T u_i \sigma_i$. In the general case, a simple calculation shows that if the spectrum of $\mathcal{V}_i^T \mathcal{U}_i$ contains m_i^+ eigenvalues equal to 1 and m_i^- equal to -1 ($m_i = m_i^+ + m_i^-$), then

$$m_i^\pm = \frac{m_i \pm \text{trace}(\mathcal{V}_i^T \mathcal{U}_i)}{2}, \quad (27)$$

i.e. the multiplicity of the eigenvalues $\pm \sigma_i$ can be easily recovered from the trace of $\mathcal{V}_i^T \mathcal{U}_i$.

To obtain the eigenvectors of A , the simplest (and more frequent) case corresponds to $m_i = 1$. In that case, the right singular vector v_i itself is an eigenvector. When some m_i is larger than 1 and $\text{trace}(\mathcal{V}_i^T \mathcal{U}_i) = m_i$ (resp. $\text{trace}(\mathcal{V}_i^T \mathcal{U}_i) = -m_i$), the m_i eigenvalues are all equal to σ_i (resp. $-\sigma_i$), and the eigenvectors are the columns of \mathcal{V}_i . In the general case $m_i > 1$, $m_i \neq m_i^\pm$, consider for each $i = 1, \dots, p$ an orthogonal diagonalization of $\mathcal{V}_i^T \mathcal{U}_i = \mathcal{W}_i J_i \mathcal{W}_i^T$, with $J_i = \text{diag}[I_{m_i^+}, -I_{m_i^-}]$ and $\mathcal{W}_i = [\mathcal{W}_i^+ | \mathcal{W}_i^-] \in \mathbb{R}^{m_i \times m_i}$ partitioned conformally to J_i . Then, denoting $\mathcal{W} = \text{diag}[\mathcal{W}_1, \dots, \mathcal{W}_p]$, one can easily check that the matrix $Q = V \mathcal{W}$ is such that

$$Q^T A Q = \text{diag}[\sigma_1 J_1, \dots, \sigma_p J_p],$$

i.e. the set of columns of the submatrix $Q_i^+ = \mathcal{V}_i \mathcal{W}_i^+ \in \mathbb{R}^{n \times m_i^+}$ (resp. $Q_i^- = \mathcal{V}_i \mathcal{W}_i^- \in \mathbb{R}^{n \times m_i^-}$) is a basis of the eigenspace corresponding to the eigenvalue σ_i (resp. $-\sigma_i$) of A . In other words, $A = Q \Lambda Q^T$ with $\Lambda = \text{diag}[\pm \sigma_i]$ is a spectral decomposition of A .

We conclude by noting that, although the right singular vectors \mathcal{V}_i have been used throughout the argument, the symmetry of A implies that similar results hold using instead the left singular vectors \mathcal{U}_i .

3.2 Clusters in finite arithmetic

We have seen in the last paragraph how to deal theoretically with groups of equal singular values. When working in finite precision, however, it is very unlikely that some of the singular values in the output of **step 2** of Algorithm SSVD come out equal. But at the same time the expected accuracy (13) determines that some of the singular values should be considered as numerically indistinguishable and treated in the spirit of § 3.1. Thus we are forced to deal with, say, k different groups Σ_i of n_i close singular values ($i = 1, \dots, k$, $n_1 + \dots + n_k = n$), which we call *clusters*⁵. The criterion to divide the singular values into clusters is crucial for the final accuracy of Algorithm SSVD. This criterion will be carefully analyzed in § 4.4, where we show

⁵For the sake of brevity, we use Σ_i to denote both the cluster of singular values and the corresponding $n_i \times n_i$ diagonal matrix.

that to achieve the accuracy (1) and (2) (see Theorems 4.3 and 4.7) it is enough to include two contiguous singular values σ_j, σ_{j+1} in the same cluster whenever

$$\frac{|\sigma_j - \sigma_{j+1}|}{\sigma_j} \leq C \kappa \epsilon, \quad (28)$$

for a suitable constant C , where

$$\kappa = \kappa(R') \max\{\kappa(X), \kappa(Y)\}$$

is the quantity (14) which came up in the error bound for the singular values computed in **step 2** of Algorithm SSVD (see § 4.4 for more on the choice of the constant C , though we advance here that in the performed numerical experiments the choice $C = 1$ gives always very satisfactory results).

For each cluster Σ_i we take matrices $U_i, V_i \in \mathbb{R}^{n \times n_i}$ whose columns are, respectively, left and right singular vectors corresponding to the singular values in Σ_i . Since the singular values in Σ_i are in general different, each of U_i and V_i is made up with several of the matrices \mathcal{U}_j and \mathcal{V}_j defined in § 3.1. Consequently, the products $\Delta_i = V_i^T U_i$ are symmetric, orthogonal and block diagonal matrices whose diagonal blocks are some of the blocks $\mathcal{V}_j^T \mathcal{U}_j$.

We conclude by noting that the number of positive, n_i^+ , and negative, n_i^- , eigenvalues with absolute values in the cluster Σ_i is still given by a formula like (27). As to the eigenvectors, things are different to § 3.1, since the diagonalization of Δ_i does not lead in general to compute eigenvectors, but just two orthonormal bases, one for the invariant subspace corresponding to the positive eigenvalues in the cluster Σ_i , and another for the negative ones. This is a fundamental issue in the error analysis for the eigenvector computations, and will be carefully explained throughout the proof of Theorem 4.4.

3.3 SYSSV Algorithm

In this section we describe Algorithm 2, the first implementation of **step 3** in Algorithm 1. The eigenvalue and the eigenvector computations are separated in the procedure in two independent parts. Doing this helps to better understand the structure of our second implementation of **step 3** in Algorithm 1, which will only insert another cluster selection routine in between the eigenvalue and the eigenvector computations. We follow LAPACK's naming convention in its name:

Algorithm 2 (SYSSV)

Input: Singular value decomposition of a symmetric matrix $A = U\Sigma V^T$.

Output: Eigenvalues $\Lambda = \text{diag}[\lambda_i]$ and eigenvectors $Q = [q_1 \dots q_n]$; $A = Q\Lambda Q^T$.

1. Decide the singular value clusters,
 $\Sigma_i = \{\sigma_{i_0}, \dots, \sigma_{i_0+n_i-1}\}, U_i, V_i, i = 1, \dots, k$, according to (28).
2. Calculate the eigenvalues using Algorithm 2.1 below.
3. Calculate the eigenvectors using Algorithm 2.2 below.

Algorithm 2.1

Input: SVD of $A = U\Sigma V^T$; Clusters $\Sigma_1, \Sigma_2, \dots, \Sigma_k$;

Output: Eigenvalues Λ ;

1. **for** each cluster, $i = 1 : k$
2. compute the diagonal elements of $\Delta_i = V_i^T U_i$
3. **if** $n_i = 1$ **then**
4. $\lambda_{i_0} = \text{sign}(\Delta_i) \sigma_{i_0}$,
5. **else**
6. **for** $j = i_0 : i_0 + n_i - 1$
7. $\lambda_j = \text{sign}[(\Delta_i)_{jj}] \sigma_j$
8. **endfor**
9. $t_i = \text{trace}(\Delta_i)$, $n_i^- = \frac{n_i - t_i}{2}$
10. **if** $\#\{(\Delta_i)_{jj} < 0\} \neq n_i^-$ **then**
11. **for** $j = i_0 : i_0 + n_i^- - 1$
12. $\lambda_j = -\sigma_j$
13. **endfor**
14. **for** $j = i_0 + n_i^- : i_0 + n_i - 1$
15. $\lambda_j = \sigma_j$
16. **endfor**
17. **endif**
18. **endfor**

Algorithm 2.2

Input: SVD of $A = U \Sigma V^T$; Clusters $\Sigma_1, \Sigma_2, \dots, \Sigma_k$; Eigenvalues Λ .

Output: Eigenvectors $Q = [q_1 \dots q_n]$

Notation:

Q_i^\pm denotes the eigenvector matrix corresponding to positive (resp. negative) e-values in Σ_i .

1. **for** each cluster, $i = 1 : k$
2. **if** $n_i = 1$ **then**
3. $q_{i_0} = v_{i_0}$
4. **else**
5. $n_i^- \equiv$ number of negative e-values in Σ_i
6. **if** $n_i^- = 0$ **then**
7. $Q_i^+ = V_i$
8. **elseif** $n_i^- = n_i$ **then**
9. $Q_i^- = V_i$
10. **else**
11. multiply $\Delta_i = V_i^T U_i$
12. diagonalize $\Delta_i = [W_i^+ W_i^-] J_i [W_i^+ W_i^-]^T$
13. $Q_i^+ = V_i W_i^+$, $Q_i^- = V_i W_i^-$
14. **endif**
15. **endif**
16. **endfor**

Some comments on this code are in order: first, we have singled out the case $n_i = 1$, although it is not needed. This is done to highlight the fact that Algorithm 2 is extremely simple in this case, all complications coming from the case $n_i > 1$.

Notice also that the code does not compute eigenvectors associated with zero eigenvalues in the case where $r = \text{rank}(A) < n$. This is due to the fact that the SVD algorithms in [7] do not compute null vectors. However, if accurate null vectors are needed, they can be obtained as the last $n - r$ columns of the orthogonal factor in a complete QR factorization of the matrix V of right singular vectors.

If large clusters are present, one can save flops in **step 11** and **step 13** of Algorithm 2.2 by employing Strassen multiplication without spoiling the accuracy of the overall algorithm. As to the diagonalization **step 12** of Algorithm 2.2, it is assumed that one performs it on a symmetrization of Δ_i . This is crucial to obtain orthonormal eigenvectors.

Notice that the eigenvalue sign assignment (**steps 6-17** of Algorithm 2.1) is done in two stages when there are clusters: first (**steps 6-8**), we assign the signs given by the diagonal elements of $\Delta_i = V_i^T U_i$, as if the singular values in Σ_i were not a cluster. If the number of assigned negative eigenvalues coincides with $n_i^- = \frac{n_i - \text{trace}(\Delta_i)}{2}$, the signs are kept. Otherwise, we proceed as described in **steps 10-17** of Algorithm 2.1. The reason for this is that the random sign assignment inside each cluster in **steps 10-17** proved to be too pessimistic in practice: although singular values inside each cluster are numerically indistinguishable according to (13), actual errors are frequently smaller than the error bounds. These smaller errors are lost if the signs of eigenvalues are randomly assigned. The modified procedure minimizes this loss of accuracy.

We finish this Section with an interesting remark on the way the signs are assigned in Algorithm 2. One might think in obtaining the sign of each eigenvalue from the Rayleigh quotients $v_i^T A v_i$, one of the most usual ways of approximating eigenvalues, instead of from $v_i^T u_i$. However, it is easy to construct examples for which the sign of $v_i^T A v_i$ is wrong, while the sign of $v_i^T u_i$ is right. We propose the following numerical example, easily reproducible in MATLAB 5.3, to the reader: generate a 100×100 symmetric Cauchy matrix with parameters $x_i = y_i \equiv r_i$, $i = 1 : 100$, where r_i is a random number chosen from a normal distribution with mean zero and variance one. Scale this matrix on both sides by the same diagonal matrix with diagonal elements $d_i = 10^{20r'_i}$, where r'_i is a random number chosen from a uniform distribution on the interval $(0.0, 1.0)$. For matrices of this kind Algorithm 3 in [6] can be used to obtain in a very simple way a RRD, $A = X D Y^T$, with forward errors fulfilling (9). Finally, applying Algorithm 3.1 of [7] to this RRD yields a SVD of A with high relative accuracy. No clusters of singular values are present in general. For several of the computed singular vectors neither $v_i^T A v_i$ nor $(v_i^T X) D (Y^T v_i)$ have the same sign of $v_i^T u_i$, which is the correct one as will be shown in Section 4 (the reader also can check this by using a symbolic package like Mathematica in very high precision). This example shows that using Rayleigh quotients may be dangerous, even in the case when the matrix is given as a RRD. Similar behaviour is not rare in other Cauchy matrices, or in random RRDs with very ill-conditioned diagonals. The use of Rayleigh quotients in the more favorable case when the matrix A is scaled in a certain particular way, is covered in [14].

4 Error analysis.

In this section we present the rounding error analysis for the eigenvalues and the eigenvectors computed by Algorithm 1. The main results in this section are the forward error bounds in Theorems 4.3 and 4.7. Both are expressed in big- O notation, without explicitly specifying the

dimensional constants involved. There are two reasons for this: first, we rely on error bounds in [7], which are written in big- O notation without explicit mention of the constants involved. Also, it is well known that the precise value of the constant is, in general, not relevant for practical purposes.

This said, the reader should be aware that in the statements of the theorems in this section we absorb moderately growing functions of the dimensions (either n , of the whole matrix, or n_i , of the clusters) as constants inside the $O(\kappa\epsilon)$. Since none of them exceeds a moderate number times n^2 , we choose not to write them explicitly in order not to complicate further the error bounds. However, the interested reader may find those corresponding to **step 3** of Algorithm 1 explicitly stated in the proofs.

The error analysis is performed in the most general case when clusters of singular values are present. This somewhat complicates the analysis, which is almost straightforward in the simple (and most likely case) of matrices whose singular values are distinct enough. The practical criterion to decide when two singular values belong to the same cluster is also discussed in detail.

This section contains the error analysis for Algorithm 1 using Algorithm 2 in **step 3**. However, it remains valid for Algorithm 1 using Algorithm 3 in **step 3**: this is trivially true for the eigenvalues, since both versions of Algorithm 1 compute the same eigenvalues. It is also true for the eigenvectors, due to the generality of the error analysis, which allows to use the new clusters appearing in Algorithm 3.

We stress that the error analysis applies *to the whole* Algorithm 1, since it relies on the backward multiplicative error formula (16), which absorbs the errors of the initial factorization in **step 1**. Although we focus on the case when the RRD is computed with the error (9), which ensures $\|E_f\| = O(\epsilon\kappa(X))$ and $\|F_f\| = O(\epsilon\kappa(Y))$, any other more general backward errors for the factorization step can be trivially incorporated into the error analysis, as explained at the end of § 2.1.

In the rest of this section we only deal with the error in nonzero eigenvalues and the corresponding eigenvectors. If the original matrix is singular, the number of zero eigenvalues is determined exactly, provided a RRD factorization fulfilling (9) is computed. As to the null vectors, it can be shown that they can be computed with error $O(\epsilon\kappa(R') \max\{\kappa(X), \kappa(Y)\})$ using the method already described after Algorithm 2.2.

We begin by fixing our model for floating point arithmetic and the notation, essential ingredients for any error analysis.

4.1 Model of arithmetic.

We use the conventional error model for floating point arithmetic:

$$\mathbf{fl}(a \odot b) = (a \odot b)(1 + \delta) \tag{29}$$

where a and b are real floating point numbers, $\odot \in \{+, -, \times, /\}$, and $|\delta| \leq \epsilon$, where ϵ is the machine precision. Moreover, we assume that neither overflow nor underflow occur. We stress that the results proved in this section still hold under a weaker error model, valid for arithmetics with no guard digit.

The error analysis below remains also valid for complex Hermitian matrices, since [17, Chapter 3] the equality (29) continues to hold for complex numbers with δ a small complex number bounded by $|\delta| = O(\epsilon)$. However, in order to simplify the presentation we only consider real symmetric matrices.

Finally, we will commit a slight abuse of notation, denoting by $\mathbf{fl}(expr)$ the computed result in finite precision of expression $expr$, instead of its rigorous meaning of the closest floating point number to $expr$.

4.2 Notation.

Letters with a hat denote computed quantities appearing in any step of Algorithm SSVD. The same letters without the hat denote their exact counterparts. It is assumed that the input of Algorithm SSVD is a real symmetric $n \times n$ matrix A , for which a RRD factorization XDY^T with small multiplicative backward error (15) can be computed.

We assume that k different clusters $\widehat{\Sigma}_i$ of n_i ($n_1 + \dots + n_k = n$) close singular values are identified through criterion (28); thus, the usual decreasing order on singular values determines the unknown exact clusters Σ_i . The singular values of one particular cluster are supposed to be different from the singular values of any other cluster. Given an index $i \in \{1, \dots, k\}$, we define

$$\Sigma_{\bar{i}} = \bigcup_{j \neq i} \Sigma_j \quad (30)$$

For each cluster Σ_i we take matrices $U_i, V_i \in \mathbb{R}^{n \times n_i}$ whose columns are, respectively, left and right singular vectors corresponding to the singular values in Σ_i . Recall that the singular values in Σ_i may be different, so both U_i and V_i will, in general, contain singular vectors corresponding to different singular values. Therefore, the remarks in §3.2 apply.

Many nontrivial choices are possible for the exact quantities U_i, V_i if A has multiple singular values in Σ_i . In that case, the results proved in this section are valid for *any* possible choice of U_i and V_i , provided their columns are singular vectors and not simply bases of the corresponding singular subspaces.

4.3 Fundamental Lemma.

The following Lemma, which is a simple consequence of the fundamental perturbation Theorem 2.3 and the multiplicative backward error formula (16) for **steps 1** and **2** of Algorithm 1, is the starting point of our error analysis.

Lemma 4.1 *Let $\widehat{U}_i, \widehat{V}_i \in \mathbb{R}^{n \times n_i}$ be the matrices of computed left and right singular vectors corresponding to the cluster of singular values $\widehat{\Sigma}_i$ computed by **steps 1–2** of Algorithm 1 applied to the symmetric matrix A . Let U_i, V_i, Σ_i be their exact counterparts. Then, there exists an exact orthogonal matrix P_i such that*

$$\sqrt{\|U_i P_i - \widehat{U}_i\|_F^2 + \|V_i P_i - \widehat{V}_i\|_F^2} \leq \frac{O(\kappa \epsilon)}{\text{relgap}(\Sigma_i, \widehat{\Sigma}_{\bar{i}})} \quad (31)$$

with κ given by (14).

Proof: Let U'_i, V'_i be the submatrices corresponding to $\widehat{\Sigma}_i$ of the exact orthogonal matrices U' and V' appearing in (16). Then, Theorem 2.3 applied to (16) guarantees that there exists an orthogonal n_i by n_i matrix P_i such that

$$\sqrt{\|U_i P_i - U'_i\|_F^2 + \|V_i P_i - V'_i\|_F^2} = \left\| \begin{bmatrix} U_i P_i - U'_i \\ V_i P_i - V'_i \end{bmatrix} \right\|_F \leq \frac{O(\kappa \epsilon)}{\text{relgap}(\Sigma_i, \widehat{\Sigma}_{\bar{i}})}.$$

Notice that

$$\sqrt{\|U_i P_i - \widehat{U}_i\|_F^2 + \|V_i P_i - \widehat{V}_i\|_F^2} = \left\| \begin{bmatrix} U_i P_i - \widehat{U}_i \\ V_i P_i - \widehat{V}_i \end{bmatrix} \right\|_F,$$

so the triangular inequality implies

$$\sqrt{\|U_i P_i - \widehat{U}_i\|_F^2 + \|V_i P_i - \widehat{V}_i\|_F^2} \leq \left\| \begin{bmatrix} U_i P_i - U_i' \\ V_i P_i - V_i' \end{bmatrix} \right\|_F + \left\| \begin{bmatrix} U_i' - \widehat{U}_i \\ V_i' - \widehat{V}_i \end{bmatrix} \right\|_F.$$

The last term in the right hand side of this inequality is $O(\epsilon)$ by (11). This concludes the proof. \blacksquare

Lemma 4.1 gives a forward error bound for simultaneous orthonormal bases of singular subspaces which depends only on the high relative accuracy algorithm used to compute the SVD (for instance Algorithm 3.1 of [7]), provided a RRD is properly computed. In other words, it only accounts for errors corresponding to **steps 1** and **2** of Algorithm 1.

The rest of the bounds obtained in this section, i.e. those corresponding to **step 3** of Algorithm 1, depend on the left hand side of (31), so it is convenient to define the quantity

$$K_i = \sqrt{\|U_i P_i - \widehat{U}_i\|_F^2 + \|V_i P_i - \widehat{V}_i\|_F^2} \quad (32)$$

for each cluster, and write all subsequent error bounds as a function of K_i . This will allow us to trace how each of the steps in Algorithm **SYSSV** contributes to the final error. From now on we assume that all quantities K_i for $i = 1, \dots, k$ are sufficiently smaller than 1, which, according to Lemma 4.1, is the case whenever the clusters of singular values are properly chosen. More precisely, all we need is that K_i be small enough to make all bounds in §4.4 and §4.5 strictly smaller than one.

4.4 Error bounds for eigenvalues and cluster criterion.

We begin by analyzing the error produced in the computation of $\text{trace}(V_i^T U_i)$ using the standard inner product algorithm.

Lemma 4.2 *Let $\widehat{U}_i, \widehat{V}_i \in \mathbb{R}^{n \times n_i}$ be the matrices of computed left and right singular vectors corresponding to the cluster of singular values $\widehat{\Sigma}_i$ computed by **steps 1-2** of Algorithm 1 applied to the symmetric matrix A . Let U_i, V_i, Σ_i be their exact counterparts. Then,*

$$\begin{aligned} \left| \mathbf{fl}(\text{trace}(\mathbf{fl}(\widehat{V}_i^T \widehat{U}_i))) - \text{trace}(V_i^T U_i) \right| &\leq \sqrt{n_i} \left(\sqrt{2} K_i + \frac{K_i^2}{2} \right) + O(\epsilon) \\ &\leq \frac{O(\kappa \epsilon)}{\text{relgap}(\Sigma_i, \widehat{\Sigma}_i)} \end{aligned} \quad (33)$$

with κ given by (14).

Proof: First observe that

$$\begin{aligned} \left| \mathbf{fl}(\text{trace}(\mathbf{fl}(\widehat{V}_i^T \widehat{U}_i))) - \text{trace}(V_i^T U_i) \right| &\leq \left| \mathbf{fl}(\text{trace}(\mathbf{fl}(\widehat{V}_i^T \widehat{U}_i))) - \text{trace}(\widehat{V}_i^T \widehat{U}_i) \right| \\ &\quad + \left| \text{trace}(\widehat{V}_i^T \widehat{U}_i) - \text{trace}(V_i^T U_i) \right|. \end{aligned} \quad (34)$$

Taking into account that the norm of the columns of \widehat{U}_i and \widehat{V}_i is close to one by (11), a straightforward error analysis [17, Chapter 3] shows that the first term in the right hand side of inequality (34) is $n_i(n + n_i)\epsilon + O(\epsilon^2)$. If P_i is the orthogonal matrix appearing in Lemma 4.1, the last term fulfils

$$\begin{aligned} \left| \text{trace}(\widehat{V}_i^T \widehat{U}_i) - \text{trace}(V_i^T U_i) \right| &= \left| \text{trace}(\widehat{V}_i^T \widehat{U}_i) - \text{trace}(P_i^T V_i^T U_i P_i) \right| \\ &\leq \sqrt{n_i} \sqrt{\sum_{k=1}^{n_i} \left| (\widehat{V}_i^T \widehat{U}_i - P_i^T V_i^T U_i P_i)_{kk} \right|^2} \\ &\leq \sqrt{n_i} \|\widehat{V}_i^T \widehat{U}_i - (V_i P_i)^T U_i P_i\|_F. \end{aligned} \quad (35)$$

Now define matrices Δ_u and Δ_v such that

$$\widehat{U}_i = U_i P_i + \Delta_u \quad \text{and} \quad \widehat{V}_i = V_i P_i + \Delta_v. \quad (36)$$

Combining (35) and (36) yields

$$\left| \text{trace}(\widehat{V}_i^T \widehat{U}_i) - \text{trace}(V_i^T U_i) \right| \leq \sqrt{n_i} (\|\Delta_u\|_F + \|\Delta_v\|_F + \|\Delta_u\|_F \|\Delta_v\|_F),$$

where we have used that $\|CD\|_F \leq \|C\|_2 \|D\|_F$ for any matrices C, D , together with the fact that the spectral norm of any matrix with orthonormal columns is one. Finally, setting $K_i = \sqrt{\|\Delta_u\|_F^2 + \|\Delta_v\|_F^2}$ as in (32), we obtain, after some direct manipulations, the desired result. ■

Notice that $\text{trace}(V_i^T U_i)$ may only take the integer values $-n_i, -n_i + 2, \dots, n_i - 4, n_i - 2, n_i$, since $V_i^T U_i$ is symmetric and orthogonal. Thus, it is sufficient that the error bound in (33) be less than one to compute *exactly* the value of $\text{trace}(V_i^T U_i)$. This can be done obtaining t_i , the nearest integer to $\mathbf{fl}(\text{trace}(\mathbf{fl}(\widehat{V}_i^T \widehat{U}_i)))$ with the parity of n_i . Then, the *integer* computation (with integer variables) of $(n_i - t_i)/2$ yields n_i^- , the *exact number of negative eigenvalues* included in the cluster Σ_i of singular values. The exact number of positive eigenvalues is obtained from the integer computation of $n_i - n_i^-$.

We stress that the conditions

$$\left| \mathbf{fl}(\text{trace}(\mathbf{fl}(\widehat{V}_i^T \widehat{U}_i))) - \text{trace}(V_i^T U_i) \right| < 1 \quad i = 1, \dots, k, \quad (37)$$

which ensure that signs are correctly assigned, determine the cluster criterion to be used in Algorithm SYSSV. Giving a rigorous criterion would require an exact knowledge of the constants involved in the big- O bound in (33), which in any case are too pessimistic in practice. Instead, we consider that the singular values in each cluster $\widehat{\Sigma}_i$ satisfy

$$\text{relgap}(\Sigma_i, \widehat{\Sigma}_i) \approx \text{relgap}(\widehat{\Sigma}_i, \widehat{\Sigma}_i) > C \epsilon \kappa(R') \max(\kappa(X), \kappa(Y))$$

for a suitable constant C . This can be obtained by imposing that two contiguous singular values $\widehat{\sigma}_j \geq \widehat{\sigma}_{j+1}$ belong to the same cluster whenever

$$\frac{|\widehat{\sigma}_j - \widehat{\sigma}_{j+1}|}{\widehat{\sigma}_j} \leq C \kappa \epsilon,$$

i.e. whenever condition (28) above holds. Choosing a large C ensures (37) and, as a consequence, that the number of positive/negative eigenvalues is correctly computed. However, a large value

for C favours the mixing of different singular values in the same cluster and, since the signs are assigned more or less randomly within each cluster, the error bound in the eigenvalues becomes roughly the product of C times the bound in the singular values. Therefore, the choice of C is subject to a certain trade-off. A sensible choice might be choosing C between 1 and 10. All the numerical experiments in section 6 have been done with $C = 1$ and the results are very satisfactory.

In any case, notice that, on one hand, the singular values are computed with the accuracy given by (16) and Theorem 2.2. On the other hand, their signs as eigenvalues of A are correctly assigned whenever the bound (33) is less than one. Hence, we have proved the main result of this subsection:

Theorem 4.3 *Let A be a $n \times n$ real symmetric matrix for which it is possible to compute a RRD fulfilling (9). Let $\lambda_1 \geq \dots \geq \lambda_n$ be the eigenvalues of A and $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_n$ be the approximations to the eigenvalues of A computed by Algorithm SSVD. Let $\hat{U}_i, \hat{V}_i \in \mathbb{R}^{n \times n_i}$ be the matrices of computed left and right singular vectors corresponding to the cluster of computed singular values $\hat{\Sigma}_i$, and let U_i, V_i, Σ_i be their exact counterparts. Assume that all clusters have been chosen according to (28), so that conditions (37) hold. Then*

$$|\lambda_j - \hat{\lambda}_j| = |\lambda_j| O(\epsilon \kappa(R') \max(\kappa(X), \kappa(Y))) \quad j = 1, \dots, n \quad (38)$$

The error bound (38) holds even for zero eigenvalues, since the *exact* number of zero eigenvalues of A is known once a RRD factorization satisfying (9) is available.

4.5 Error bounds for eigenvectors.

In this section we obtain bounds on the distance between *bases* of invariant subspaces. Although it is more usual to bound the sines of the canonical angles between the exact and the computed invariant subspaces [26], we choose to compare the bases themselves because, as explained before Theorem 2.3, bases play an essential role both in Algorithm SYSSV and in its error analysis. However, usual $\sin \Theta$ bounds easily follow from Theorem 4.7, since distances between bases and canonical angles between subspaces are closely related [26, Theorem I.5.2, Theorem II.4.11] and the same bounds hold for both, up to a factor $\sqrt{2}$ in Frobenius norm.

One important issue in the subsequent analysis comes from **step 12** of Algorithm 2.2 in which the $n_i \times n_i$ matrix $\hat{V}_i^T \hat{U}_i$ is orthogonally diagonalized for each cluster $\hat{\Sigma}_i$. Lemma 4.1 shows that the matrices \hat{U}_i, \hat{V}_i of computed singular vectors are not reliable approximations of the matrices of exact singular vectors U_i, V_i , but just reliable approximations of $U_i P_i$ and $V_i P_i$, with P_i the unknown $n_i \times n_i$ orthogonal matrix in Lemma 4.1. Hence, we are forced in practice to diagonalize approximations to matrices $P_i^T V_i^T U_i P_i$. Next Theorem in exact arithmetic shows that this is enough to get orthonormal bases of invariant subspaces, although not for obtaining eigenvectors.

Theorem 4.4 *Let A be a symmetric $n \times n$ matrix and $U_i, V_i \in \mathbb{R}^{n \times n_i}$ be matrices of left and right singular vectors of A corresponding to a cluster of nonzero singular values Σ_i , different from the rest of the singular values of A . Let P_i be any $n_i \times n_i$ orthogonal matrix, and consider any orthogonal diagonalization of the $n_i \times n_i$ orthogonal and symmetric matrix $P_i^T V_i^T U_i P_i$ partitioned as*

$$P_i^T V_i^T U_i P_i = [W_i^+ W_i^-] \begin{bmatrix} I_{n_i^+} & 0 \\ 0 & -I_{n_i^-} \end{bmatrix} [W_i^+ W_i^-]^T, \quad (39)$$

where I_s denotes the $s \times s$ identity matrix and $n_i^+ + n_i^- = n_i$. Then the columns of $V_i P_i W_i^+$ (resp. $V_i P_i W_i^-$) form an orthonormal basis of the invariant subspace of A corresponding to the positive (resp. negative) eigenvalues whose absolute values are in Σ_i .

Proof: Without loss of generality we may consider the SVD of A partitioned in only two blocks,

$$A = [U_1 U_2] \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix} [V_1 V_2]^T, \quad (40)$$

where no special order is assumed on the singular values. Here Σ_1 corresponds to the cluster Σ_i to be studied and Σ_2 to the remaining ones $\Sigma_{\bar{i}}$ defined as in (30). The matrix P_i will be denoted just by P , and the matrices W_i^\pm in (39) by W_\pm .

As mentioned in § 3.2, $V_1^T U_1$ is orthogonal, symmetric and block diagonal with the size of the blocks fixed by the groups of equal singular values inside Σ_1 . The matrix $V_1^T U_1 \Sigma_1$ is also symmetric with the same block diagonal structure of $V_1^T U_1$. An orthogonal diagonalization for each block of $V_1^T U_1$ leads to an orthogonal diagonalization of the full matrix $V_1^T U_1$ with eigenvectors which are also eigenvectors of $V_1^T U_1 \Sigma_1$. In this situation, the eigenvectors of $V_1^T U_1$ corresponding to the eigenvalue 1 (resp. -1) are the eigenvectors of $V_1^T U_1 \Sigma_1$ corresponding to positive (resp. negative) eigenvalues with absolute values in Σ_1 . From this we deduce that the invariant subspaces corresponding to positive (resp. negative) eigenvalues of matrices $P^T V_1^T U_1 P$ and $P^T V_1^T U_1 \Sigma_1 P$ coincide. Once this is taken into account, the rest of the proof reduces to some easy block manipulations.

Combining (40) and $V_2^T U_1 = 0$ from (24), we obtain

$$A V_1 P = U_1 \Sigma_1 P = [V_1 V_2] \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix} U_1 \Sigma_1 P = V_1 P (P^T V_1^T U_1 \Sigma_1 P). \quad (41)$$

Splitting the spectrum into positive and negative eigenvalues, we orthogonally diagonalize

$$P^T V_1^T U_1 \Sigma_1 P = [Q_+ Q_-] \begin{bmatrix} D_+ & 0 \\ 0 & D_- \end{bmatrix} [Q_+ Q_-]^T,$$

and from (41) we obtain

$$A(V_1 P Q_+) = (V_1 P Q_+) D_+ \quad \text{and} \quad A(V_1 P Q_-) = (V_1 P Q_-) D_-. \quad (42)$$

Now, we know that $\text{col}(Q_\pm) = \text{col}(W_\pm)$, and since the columns of Q_\pm and W_\pm are orthonormal, there exist square orthogonal matrices T_\pm such that $W_\pm = Q_\pm T_\pm$. Combining this and (42) we obtain

$$A(V_1 P W_\pm) = (V_1 P W_\pm) (T_\pm^T D_\pm T_\pm),$$

which proves the Theorem. ■

Once the previous Theorem has been proved, the rest of the section is organized in the following three steps:

1. Although Lemma 4.1 guarantees that \widehat{U}_i and \widehat{V}_i are close to $U_i P_i$ and $V_i P_i$, provided the clusters have been properly chosen, this does not mean that $\widehat{\Delta}_i = \mathbf{fl}(\widehat{V}_i^T \widehat{U}_i)$ in **step 11** of Algorithm 2.2 is symmetric. Let \widehat{S}_i be the symmetric matrix obtained by replacing the upper triangular part of $\widehat{\Delta}_i$ by its lower triangular part. Lemma 4.5 bounds the difference

between \widehat{S}_i and the exact symmetric matrix $P_i^T V_i^T U_i P_i$. Notice that if any driver routine of LAPACK [1] for the symmetric eigenvalue problem is used in **step 12** of Algorithm 2.2, just the upper (or lower) triangular part of $\widehat{\Delta}_i$ is stored. Hence, the symmetrization step does not require any additional work.

2. Lemma 4.6 relates the computed orthogonal eigendecomposition of \widehat{S}_i with an exact eigendecomposition of $P_i^T V_i^T U_i P_i$. It is shown that exact matrices W_i^\pm in (39) can be chosen close enough to the corresponding computed matrices \widehat{W}_i^\pm in **step 12** of Algorithm 2.2.
3. Finally, the main Theorem 4.7 bounds the difference between the $n \times n_i^\pm$ matrices $\mathbf{fl}(\widehat{V}_i \widehat{W}_i^\pm)$ computed in **step 13** of Algorithm 2.2 and some orthonormal bases of exact invariant subspaces of A . This result is a simple consequence of Lemma 4.1 and Lemma 4.6.

The bottom line after these three steps is that **step 3** of Algorithm 1 produces errors of the order of K_i , the quantity defined in (32), whose upper bound (31) depends only on the errors in **steps 1** and **2** of Algorithm 1.

Lemma 4.5 *Let $\widehat{U}_i, \widehat{V}_i \in \mathbb{R}^{n \times n_i}$ be the matrices of computed left and right singular vectors corresponding to the cluster of singular values $\widehat{\Sigma}_i$ computed by **step 1-2** of Algorithm 1 applied to the symmetric matrix A . Let U_i, V_i, Σ_i be their exact counterparts. Let \widehat{S}_i be a symmetrization of the floating point matrix $\widehat{\Delta}_i = \mathbf{fl}(\widehat{V}_i^T \widehat{U}_i)$ obtained by replacing the upper triangular part of $\widehat{\Delta}_i$ by its lower triangular part, or viceversa. Then an orthogonal $n_i \times n_i$ matrix P_i exists such that*

$$\begin{aligned} \|\widehat{S}_i - P_i^T V_i^T U_i P_i\|_F &\leq 2K_i + \frac{K_i^2}{\sqrt{2}} + O(\epsilon) \\ &\leq \frac{O(\epsilon \kappa(R') \max(\kappa(X), \kappa(Y)))}{\text{relgap}(\Sigma_i, \widehat{\Sigma}_i)} \end{aligned} \quad (43)$$

Proof: First observe that

$$\|\mathbf{fl}(\widehat{V}_i^T \widehat{U}_i) - P_i^T V_i^T U_i P_i\|_F \leq \|\mathbf{fl}(\widehat{V}_i^T \widehat{U}_i) - \widehat{V}_i^T \widehat{U}_i\|_F + \|\widehat{V}_i^T \widehat{U}_i - P_i^T V_i^T U_i P_i\|_F$$

where P_i is the orthogonal matrix appearing in Lemma 4.1. Standard error analysis of usual matrix multiplication [17], and the fact that the columns of \widehat{U}_i and \widehat{V}_i are almost orthonormal by (11), shows that the first term of the right hand side of the previous inequality is bounded by $n n_i \epsilon + O(\epsilon^2)$. The last term can be bounded as in the proof of Lemma 4.2, so we obtain

$$\|\mathbf{fl}(\widehat{V}_i^T \widehat{U}_i) - P_i^T V_i^T U_i P_i\|_F \leq (\sqrt{2}K_i + \frac{K_i^2}{2}) + O(\epsilon).$$

We write $\mathbf{fl}(\widehat{V}_i^T \widehat{U}_i) = \widehat{L} + \widehat{D} + \widehat{R}$ as the sum of its strict lower triangular part, its diagonal part and its strict upper triangular part. The same is done for the symmetric matrix $P_i^T V_i^T U_i P_i = L + D + L^T$, so the previous equation yields

$$\sqrt{\|(\widehat{L} + \widehat{D}) - (L + D)\|_F^2 + \|\widehat{R} - L^T\|_F^2} \leq (\sqrt{2}K_i + \frac{K_i^2}{2}) + O(\epsilon). \quad (44)$$

The same inequality holds for $\sqrt{\|\widehat{L} - L\|_F^2 + \|\widehat{D} + \widehat{R} - (D + L^T)\|_F^2}$. On the other hand

$$\|\widehat{S}_i - P_i^T V_i^T U_i P_i\|_F = \sqrt{\|(\widehat{L} + \widehat{D}) - (L + D)\|_F^2 + \|\widehat{L}^T - L^T\|_F^2}.$$

Combining this equation with (44) proves the Lemma. ■

Errors in the diagonalization **step 12** of Algorithm 2.2 are now analyzed. Notation and definitions of previous Lemma are used.

Lemma 4.6 *Let $\widehat{W}_i \widehat{\Lambda}_i \widehat{W}_i^T$ be the computed orthogonal spectral decomposition of the symmetric $n_i \times n_i$ matrix \widehat{S}_i using any LAPACK subroutine for the symmetric eigenproblem [1, § 2.3.4.1]. Then, there exists a matrix E_i , an orthogonal matrix Z_i and an orthogonal matrix P_i such that*

$$P_i^T V_i^T U_i P_i + E_i = Z_i \widehat{\Lambda}_i Z_i^T \quad (45)$$

where

$$\|Z_i - \widehat{W}_i\|_2 \leq O(\epsilon) \quad \text{and} \quad \|E_i\|_F \leq 2K_i + \frac{K_i^2}{\sqrt{2}} + O(\epsilon). \quad (46)$$

Moreover, if \widehat{W}_i^+ (resp. \widehat{W}_i^-) is the submatrix of \widehat{W}_i with columns corresponding to the positive (resp. negative) elements of $\widehat{\Lambda}_i$, then there exist matrices W_i^+ , W_i^- fulfilling (39) such that

$$\begin{aligned} \|\widehat{W}_i^\pm - W_i^\pm\|_F &\leq 2\sqrt{2}K_i + K_i^2 + O(\epsilon) \\ &= \frac{O(\epsilon \kappa(R') \max(\kappa(X), \kappa(Y)))}{\text{relgap}(\Sigma_i, \widehat{\Sigma}_i)} \end{aligned} \quad (47)$$

Proof: Using the results in [1, § 4.7.1] we see that there exist an orthogonal matrix Z_i and a matrix E'_i such that

$$\widehat{S}_i + E'_i = Z_i \widehat{\Lambda}_i Z_i^T \quad (48)$$

where

$$\|Z_i - \widehat{W}_i\|_2 \leq O(\epsilon) \quad \text{and} \quad \|E'_i\|_2 \leq O(\epsilon) \|\widehat{S}_i\|_2.$$

Let P_i be the orthogonal matrix appearing in Lemmas 4.1 and 4.5. The spectral norm of the orthogonal matrix $P_i^T V_i^T U_i P_i$ is equal to one, so (43) implies $\|\widehat{S}_i\|_2 = 1 + \beta$, with $|\beta| \leq 2K_i + K_i^2/\sqrt{2} + O(\epsilon)$. Thus $\|E'_i\|_2 = O(\epsilon)$. Now, expressions (45) and (46) are easily proved using Lemma 4.5, noting by (48) that

$$P_i^T V_i^T U_i P_i + \widehat{S}_i - P_i^T V_i^T U_i P_i + E'_i = Z_i \widehat{\Lambda}_i Z_i^T,$$

and defining

$$E_i = \widehat{S}_i - P_i^T V_i^T U_i P_i + E'_i.$$

We finally prove (47): let W_i^\pm be matrices fulfilling (39) and Z_i^+ (resp. Z_i^-) be a submatrix of Z_i corresponding to the positive (resp. negative) elements of $\widehat{\Lambda}_i$. We assume that K_i is small enough to imply $\|E_i\|_2 < 1$, so the eigenvalues equal to 1 (resp. -1) of $P_i^T V_i^T U_i P_i$ remain positive (resp. negative) in $\widehat{\Lambda}_i$. This can be seen by applying Weyl's eigenvalue perturbation theorem to (45) (see for instance [26, Corollary IV.4.10]). Thus, Davis and Kahan's $\sin \Theta$ theorem for variations of invariant subspaces of Hermitian matrices [5] applied to (45) leads to

$$\|\sin \Theta(W_i^+, Z_i^+)\|_F \leq \frac{\|E_i\|_F}{\min_{\substack{\mu < 0 \\ \mu \in \widehat{\Lambda}_i}} |1 - \mu|} \leq \|E_i\|_F, \quad (49)$$

where the matrix $\Theta(W_i^+, Z_i^+)$ is the matrix of the canonical angles between the column space of W_i^+ and the column space of Z_i^+ . Theorem II.4.11 in [26], (49) and (46) show that it is possible to choose W_i^+ such that

$$\begin{aligned} \|W_i^+ - Z_i^+\|_F &= \sqrt{\|\sin \Theta(W_i^+, Z_i^+)\|_F^2 + \|I - \cos \Theta(W_i^+, Z_i^+)\|_F^2} \\ &\leq \sqrt{2} \|\sin \Theta(W_i^+, Z_i^+)\|_F \\ &\leq \sqrt{2} \|E_i\|_F \\ &\leq 2\sqrt{2}K_i + K_i^2 + O(\epsilon). \end{aligned} \quad (50)$$

Similar results hold for W_i^- and Z_i^- . We finish the proof by noting that

$$\|\widehat{W}_i^\pm - W_i^\pm\|_F \leq \|\widehat{W}_i^\pm - Z_i^\pm\|_F + \|Z_i^\pm - W_i^\pm\|_F.$$

The first term of the right hand side is $O(\epsilon)$ by (46), and the second one is bounded in (50). ■

We conclude with the main result on rounding errors for eigenvectors computed in **step 13** of Algorithm 2.2. Previous notation and definitions are used.

Theorem 4.7 *Let A be a $n \times n$ real symmetric matrix of rank r for which it is possible to compute a RRD fulfilling (9). Let $\widehat{\Sigma}_i$ be a cluster of nonzero computed singular values of A using **steps 1-2** of Algorithm 1 and Σ_i be the corresponding cluster of exact singular values. Then there exist matrices Q_i^+ and Q_i^- , whose columns form orthonormal bases of the invariant subspaces of A corresponding, respectively, to the positive and negative eigenvalues of A with absolute values in Σ_i , such that*

$$\begin{aligned} \|\mathbf{f1}(\widehat{V}_i \widehat{W}_i^+) - Q_i^+\|_F &\leq (2\sqrt{2} + 1)(K_i + K_i^2) + K_i^3 + O(\epsilon) \\ &= \frac{O(\epsilon \kappa(R') \max(\kappa(X), \kappa(Y)))}{\text{relgap}(\Sigma_i, \widehat{\Sigma}_i)}, \end{aligned} \quad (51)$$

with and equal bound for $\|\mathbf{f1}(\widehat{V}_i \widehat{W}_i^-) - Q_i^-\|_F$.

Moreover, let $\widehat{Q} = [\mathbf{f1}(\widehat{V}_1 \widehat{W}_1^+) \mathbf{f1}(\widehat{V}_1 \widehat{W}_1^-) \dots \mathbf{f1}(\widehat{V}_k \widehat{W}_k^+) \mathbf{f1}(\widehat{V}_k \widehat{W}_k^-)]$ be the $n \times r$ matrix whose columns are the bases of all considered invariant subspaces of A computed using Algorithm 1. Then there exists an $n \times r$ matrix B with exact orthonormal columns such that

$$\|\widehat{Q} - B\|_F = O(\epsilon). \quad (52)$$

Proof: Let \widehat{V}_i be the matrix of computed right singular vectors corresponding to the cluster $\widehat{\Sigma}_i$, and V_i its exact counterpart. Let W_i^\pm , \widehat{W}_i^\pm and P_i be the matrices appearing in Lemmas 4.6 and 4.1. By Theorem 4.4, the columns of $Q_i^+ \equiv V_i P_i W_i^+$ and $Q_i^- \equiv V_i P_i W_i^-$ are orthonormal bases of the invariant subspaces of A corresponding, respectively, to the positive and negative eigenvalues of A with absolute values in Σ_i .

Note also that

$$\|\mathbf{f1}(\widehat{V}_i \widehat{W}_i^\pm) - V_i P_i W_i^\pm\|_F \leq \|\mathbf{f1}(\widehat{V}_i \widehat{W}_i^\pm) - \widehat{V}_i \widehat{W}_i^\pm\|_F + \|\widehat{V}_i \widehat{W}_i^\pm - V_i P_i W_i^\pm\|_F. \quad (53)$$

The first term of the right hand side is bounded by $n_i \sqrt{n_i n_i^\pm} \epsilon + O(\epsilon^2)$ using the standard error analysis of usual matrix multiplication [17] and the fact that the columns of \widehat{V}_i and \widehat{W}_i^\pm

are almost orthonormal by (11) and (46). For the second term we proceed as follows: define matrices Δ_v and Δ_w^\pm by

$$\widehat{V}_i = V_i P_i + \Delta_v \quad \text{and} \quad \widehat{W}_i^\pm = W_i^\pm + \Delta_w^\pm,$$

where $\|\Delta_v\|_F \leq K_i$ by (32) and $\|\Delta_w^\pm\|_F \leq 2\sqrt{2}K_i + K_i^2 + O(\epsilon)$ by (47). Thus

$$\begin{aligned} \|\widehat{V}_i \widehat{W}_i^\pm - V_i P_i W_i^\pm\|_F &\leq \|\Delta_v\|_F + \|\Delta_w^\pm\|_F + \|\Delta_v\|_F \|\Delta_w^\pm\|_F \\ &\leq (2\sqrt{2} + 1)(K_i + K_i^2) + K_i^3 + O(\epsilon). \end{aligned}$$

Combining this with (53) proves (51).

Finally (52) follows from the well-known fact that finite precision matrix multiplication of matrices with columns orthonormal up to $O(\epsilon)$ yields a matrix with columns orthonormal up to $O(\epsilon)$. \blacksquare

As announced in the Introduction, the eigenvector error bounds we derive suffer from an important drawback: they depend on *relgap* (22) between singular values which is less than or equal to the natural relative gap between eigenvalues, the one expected for the symmetric eigenproblem. This is an unavoidable consequence of the non-symmetric character of Algorithm SSVD. This drawback, however, can be partially solved applying Theorem 4.7 to certain new singular value clusters chosen as described in section 5.

It is worth observing that Theorem 4.7 does not guarantee that the columns of the matrices $\mathbf{fl}(\widehat{V}_i \widehat{W}_i^\pm)$ computed by Algorithm SSVD always approximate *eigenvectors* of A . This can only be ensured in three cases: when there is no cluster ($n_i = 1$), when all eigenvalues in the cluster have the same sign, or when the cluster contains eigenvalues of both signs with either $n_i^+ = 1$ or $n_i^- = 1$. In this last case, either $\mathbf{fl}(\widehat{V}_i \widehat{W}_i^+)$ or $\mathbf{fl}(\widehat{V}_i \widehat{W}_i^-)$ approximates an eigenvector of A . In any other situation, the columns of $\mathbf{fl}(\widehat{V}_i \widehat{W}_i^\pm)$ do not approximate eigenvectors, but just orthonormal bases of the invariant subspaces of A corresponding to either the positive or the negative eigenvalues with absolute values in the cluster. However, provided the clusters of singular values are chosen according to criterion (28), this does not represent any drawback, because the eigenvectors in the corresponding invariant subspaces are computed by any symmetric eigensolver (including the J-orthogonal algorithm [27, 23]) with large errors due to the presence of very small relative gaps between the eigenvalues inside the clusters. No need to say that the J-orthogonal algorithm also computes accurate bases of invariant subspaces, due to its backward stability.

We conclude with an interesting remark concerning the discussion in the previous paragraph. Consider, for simplicity, that according to criterion (28) a cluster of two singular values, one corresponding to a positive eigenvalue and the other to a negative one, has been found. Then the bound in Theorem 4.7 implies that Algorithm SSVD computes *both* eigenvectors with an error governed by the relative gap between the cluster and the singular values outside the cluster. This can be much larger than the relative gap between the singular values inside the cluster. Thus, the presence of clusters reduces the errors in the computed eigenvectors. We will take more advantage of this property in Section 5.

5 Computing more accurate eigenvectors

The error in the eigenvectors computed by Algorithm 2 (SYSSV) is governed by the singular value relative gap (see Theorem 4.7), which is less than or equal to the natural eigenvalue relative gap.

In this section we propose an alternative algorithm to SYSSV which computes eigenvectors with the error (10) announced in the Introduction. As we will see, the underlying idea is very simple and does not require a new error analysis, but just taking advantage of the generality of the error analysis in section 4. Let us remark that the eigenvalue computation (**Steps 1-2** in Algorithm 2) will stay the same, in this section we will only modify the computation of the eigenvectors. The general case, when clusters of singular values of arbitrary dimension are present, will be considered.

First, note that in Algorithm 2 the eigenvalues are computed before computing the eigenvectors. The relative error in the eigenvalues is of order $O(\epsilon\kappa(R') \max(\kappa(X), \kappa(Y)))$ provided the clusters are chosen according to criterion (28). A second important remark is that the error analysis performed in section 4 for the eigenvectors is independent of the error analysis for the eigenvalues, both being valid under the hypothesis that the quantities K_i defined in (32) are sufficiently smaller than 1. As Lemma 4.1 shows this is achieved by defining clusters which yield large enough $relgap(\Sigma_i, \widehat{\Sigma}_i)$, but whenever this condition is fulfilled different clusters, i.e. different K_i , can be chosen to compute the eigenvectors using Algorithm 2.2. Theorem 4.7 still applies and will provide a smaller error bound whenever the new clusters *for the eigenvector computation* have larger $relgaps$ than the ones chosen according to (28).

The underlying idea is very simple: Let Σ_i be one of the singular value cluster chosen according to (28), and let Λ_i^+ (resp. Λ_i^-) be the corresponding clusters of positive (resp. negative) eigenvalues with absolute values in Σ_i . Then $relgap(\Sigma_i, \widehat{\Sigma}_i)$ can be much worse than the minimum of the two eigenvalue relative gaps associated to Σ_i only in the case in which Σ_i is signed (all the eigenvalues of the same sign), and the closest (in the relative sense) cluster to Σ_i , let us say $\Sigma_{cl(i)}$, is oppositely signed. With no loss of generality it can be supposed that $\Sigma_i = \Lambda_i^+$, therefore $\Sigma_{cl(i)} = -\Lambda_{cl(i)}^-$. If Σ_i and $\Sigma_{cl(i)}$ are joined to form a new cluster $\{\Lambda_i^+, -\Lambda_{cl(i)}^-\}$ with a larger $relgap$, the bound (51) will improve **separately** for the bases of *exactly the same two invariant subspaces* associated with Λ_i^+ and $\Lambda_{cl(i)}^-$, calculated by Algorithm 2.2 applied to the new set of clusters. Therefore, nothing is lost by merging clusters of this kind and the error bound (51) can improve by *joining* close adjacent clusters in such a way that $relgap$ increases.

The whole idea can be summarized in the following algorithm, alternative to SYSSV:

Algorithm 3 (SYSSVR)

Input: Singular value decomposition of a symmetric matrix $A = U\Sigma V^T$.

Output: Eigenvalues $\Lambda = diag[\lambda_i]$ and eigenvectors $Q = [q_1 \dots q_n]$; $A = Q\Lambda Q^T$.

1. Decide the singular value clusters, $\{\Sigma_i, U_i, V_i\}_{i=1}^k$, according to (28).
2. Calculate the eigenvalues using Algorithm 2.1.
3. Use Algorithm 3.1 in §5.2 to merge, when necessary, some pairs of clusters to form a new set $\{\Sigma_i, U_i, V_i\}_{i=1}^q$ of clusters, according to the strategy developed in this section.
4. Calculate the eigenvectors using Algorithm 2.2 on the new set of clusters.

We will show in this section that Algorithm 3 (SYSSVR) guarantees that the error in the computed basis of the invariant subspace corresponding to each cluster of eigenvalues $\widehat{\Lambda}_i$ of the symmetric matrix A is smaller than

$$\frac{O(\epsilon\kappa(R') \max(\kappa(X), \kappa(Y)))}{\min\{relgap(\widehat{\Lambda}_i), relgap(\widehat{\Lambda}_{cl(i)})\}}, \quad (54)$$

where the eigenvalue relative gap in the denominator corresponds, either to the cluster $\widehat{\Lambda}_i$ under consideration, or to the cluster $\widehat{\Lambda}_{cl(i)}$ whose eigenvalues have different sign, but are the closest (in relative sense) in absolute value. This result will be proved in Theorem 5.12 and generalizes to invariant subspaces the error bound (10) appearing in the introduction for eigenvectors.

The rest of this section is organized as follows: some relationships between eigenvalue and singular value relative gaps are proved in §5.1. This is necessary if (54) has to be proved using Theorem 4.7, which only deals with singular value relative gaps. First we show in Theorem 5.5 that in the case of an *unsigned* cluster (a cluster containing singular values corresponding to positive and negative eigenvalues), the singular value relative gap of the cluster is not worse, up to a moderate constant, than an eigenvalue relative gap. Theorem 5.6 proves that this also happens to the relative gap of a signed cluster if the closest cluster is not signed of the opposite sign. Thus for clusters of these two kinds (54) holds, and it is not necessary to join them to any other cluster.

In the rest of subsection 5.1 we will study the case of a signed cluster whose closest cluster is oppositely signed. In all the theorems it will be supposed that the singular value relative gap is sufficiently smaller than the eigenvalue relative gap, otherwise it is trivial that (54) is reached. With these assumptions (54) is always achieved, either by joining clusters if the singular value relative gap improves (Theorem 5.7), or if not, by doing nothing (Theorem 5.9). Finally Theorem 5.10 proves that it is not necessary to join more than two clusters. Let us remark that only in the case of the hypotheses of Theorem 5.7 something different than Algorithm 2 has to be done to get (54).

With these results at hand, in subsection 5.2 we will implement a routine (Algorithm 3.1) that joins singular value clusters, previously chosen according to (28), following the strategy of the theorems in subsection 5.1. Then to these new clusters Algorithm 2.2 is applied and Theorem 5.12 proves that (54) will be achieved for the computed bases of the invariant subspaces.

Here, as in Section 4 only clusters of nonzero singular values will be considered. Apart from the reasons stated in Section 4, it should be remarked that a cluster of zero singular values is at the same time a cluster of zero eigenvalues, and both its eigenvalue and singular value relative gaps are equal to 1. Thus for such cluster an error bound $O(\epsilon\kappa(R')\max(\kappa(X),\kappa(Y)))$ holds, which is better than (54). Moreover a cluster of zero singular values is as far as possible, in relative distance, from any other cluster, thus joining it to other cluster makes no sense.

5.1 Eigenvalue versus singular value relative gaps.

Throughout this section we consider a set of real numbers $\Lambda = \{\lambda_1, \dots, \lambda_n\}$ decreasingly ordered, i.e., $\lambda_1 \geq \dots \geq \lambda_n$, and the set of their moduli, $\Sigma = \{\sigma_1, \dots, \sigma_n\}$ also in decreasing order, i.e., $\sigma_1 \geq \dots \geq \sigma_n \geq 0$. Let Π be the index permutation such that $\sigma_i = |\lambda_{\Pi(i)}|$. Whenever we consider a subset $\Sigma_1 = \{\sigma_{i+1}, \sigma_{i+2}, \dots, \sigma_{i+d_1}\}$ of Σ we will denote by $\Lambda_1 = \{\lambda_{\Pi(i+1)}, \dots, \lambda_{\Pi(i+d_1)}\}$ the corresponding subset of Λ ; moreover we will call Λ_1^+ (resp. Λ_1^-) the set of positive (resp. negative) elements of Λ_1 . It is worth to think in Λ and Σ as being, respectively, the set of eigenvalues and singular values of the real symmetric matrix A studied in Section 4, but notice that the following results are proved using only elementary properties of real numbers without any reference to spectral properties. Thus, the proofs of the theorems appearing in this subsection are all of them elementary but sometimes long and involved, mainly due to dealing with clusters containing more than one element. By this reason, we postpone the proofs to Appendix B, and only those of the more intricate results, Theorems 5.7 and 5.9, will be explained, the rest are similar.

Our definitions of relative gaps, see (4) and (8), are convenient and appealing in numerical

analysis but the lack of symmetry in relative errors of the type $|\sigma_j - \sigma_i|/\sigma_i$ is unpleasant from a mathematical point of view and complicates somewhat the statement of the results (see more on these questions and definitions of true relative mathematical distances in [19, 20]). In this sense, an effort to state the theorems with the idea of being applied to develop Algorithm 3.1 has been done.

We begin by a general definition of cluster:

Definition 5.1 Let C_l be a real number such that $0 \leq C_l < 1$. The subset $\Sigma_1 = \{\sigma_{i+1}, \sigma_{i+2}, \dots, \sigma_{i+d_1}\}$ of Σ is called a cluster of tolerance C_l if

1. $(\sigma_j - \sigma_{j+1}) \leq C_l \sigma_j$ for $j = i + 1, \dots, i + d_1 - 1$,
2. $(\sigma_i - \sigma_{i+1}) > C_l \sigma_i$ and $(\sigma_{i+d_1} - \sigma_{i+d_1+1}) > C_l \sigma_{i+d_1}$, whenever all the labels belong to $\{1, 2, \dots, n\}$, otherwise the corresponding inequality does not appear in the definition.

Notice that in case of a cluster of dimension 1 ($d_1 = 1$) the first condition is empty. Notice also that this definition includes the clusters of singular values chosen in Algorithm 2 SYSSV, according to criterion (28), for $C_l = \epsilon \kappa(R') \max\{\kappa(X), \kappa(Y)\}$. The condition $C_l < 1$ appearing in Definition 5.1 is necessary, otherwise the whole set Σ would be always a trivial cluster, independently of the distribution of its elements.

Now we define relative gaps for subsets of Λ and Σ . For the sake of simplicity we will only use one argument.

Definition 5.2 Let Λ_2 and Σ_1 be any subsets of, respectively, Λ and Σ . We define the following relative gaps for both subsets:

1.
$$rg(\Lambda_2) = \min_{\substack{\lambda_k \in \Lambda_2 \\ \lambda_q \notin \Lambda_2}} \frac{|\lambda_q - \lambda_k|}{|\lambda_k|}.$$

2.
$$relgap(\Lambda_2) = \min\{rg(\Lambda_2), 1\}.$$

3.
$$rg(\Sigma_1) = \min_{\substack{\sigma_k \in \Sigma_1 \\ \sigma_q \notin \Sigma_1}} \frac{|\sigma_q - \sigma_k|}{\sigma_k}.$$

4.
$$relgap(\Sigma_1) = \min\{rg(\Sigma_1), 1\}.$$

Given a subset Σ_1 of Σ , the relationship between the $relgap(\Sigma_1)$ appearing in the definition 5.2 and $relgap$ as defined by (22) and (20) is

$$relgap(\Sigma_1) = relgap(\Sigma_{\bar{1}}, \Sigma_1), \tag{55}$$

where the notation introduced in (30) has been used. Similar comments apply to rg defined in (20) and rg defined above. Although $relgap(\Sigma_1, \Sigma_{\bar{1}})$ is the relative gap appearing in the error analysis of Section 4, we have found it simpler, both from a theoretical and computational point of view, to deal with $relgap(\Sigma_i)$, which has the elements of the cluster being analyzed in the

denominators of the relative errors⁶. Both choices are equivalent, as shown in (23), and, on the other hand, it is possible to reformulate Theorem 2.3 using $relgap(\Sigma_i)$.

The error bounds for invariant subspaces computed using the J-orthogonal and SSVD algorithms are controlled by the relative gaps $relgap$, of eigenvalues and singular values respectively, in the previous definition (see Theorem 4.7 and [23, p.7]). However in the following it is simpler and more general to use the relative gaps rg . At the end of this section it will be shown that theorems obtained for rg imply easily results for $relgap$.

We start with this simple Lemma:

Lemma 5.3 *Let $\Sigma_1 = \{\sigma_{i+1}, \sigma_{i+2}, \dots, \sigma_{i+d_1}\}$ be a subset of consecutive elements of Σ . Then*

$$rg(\Sigma_1) = \min \left\{ \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}}, \frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1}} \right\},$$

where if the index i or $i + d_1 + 1$ does not belong to $\{1, \dots, n\}$ the corresponding term does not appear in the minimum.

This lemma allows a natural definition of the *closest cluster to Σ_1 in the relative sense*

Definition 5.4 *Let $\Sigma_1 = \{\sigma_{i+1}, \sigma_{i+2}, \dots, \sigma_{i+d_1}\}$ be a cluster of tolerance C_l , we define its relative closest cluster $\Sigma_{cl(1)}$ as the cluster of tolerance C_l containing σ_i if $rg(\Sigma_1) = (\sigma_i - \sigma_{i+1})/\sigma_{i+1}$, or the one containing σ_{i+d_1+1} if $rg(\Sigma_1) = (\sigma_{i+d_1} - \sigma_{i+d_1+1})/\sigma_{i+d_1}$.*

It is seen from Lemma 5.3 that, with the possible exception of the cluster containing the smallest singular value, $rg(\Sigma_1) \leq 1$ and then $rg(\Sigma_1) = relgap(\Sigma_1)$. Obviously the last equality also holds whenever $rg(\Sigma_1) < 1$, a condition appearing frequently in the results of this section.

Our first result deals with the case of clusters containing singular values corresponding to positive and negative eigenvalues. This theorem shows that in this case the singular value relative gap of the cluster is not worse, up to a moderate constant, than an eigenvalue relative gap. Thus for clusters of singular values of this kind (54) holds, and it is not necessary to join them to any other cluster.

Theorem 5.5 *Let Σ_1 be a cluster of singular values of tolerance C_l with d_1 elements such that $(d_1 - 1)C_l < 1$, and assume that Λ_1 contains both positive and negative elements. Then*

$$\min\{rg(\Lambda_1^+), rg(\Lambda_1^-)\} \leq \frac{1}{1 - (d_1 - 1)C_l} \left(1 + \frac{(d_1 - 1)C_l}{rg(\Sigma_1)} \right) rg(\Sigma_1).$$

First, notice that Theorem 5.5 greatly simplifies in the case of one-dimensional ($d_1 = 1$) clusters. However some remarks about the bound in previous theorem are in order when $d_1 > 1$. The assumption $(d_1 - 1)C_l < 1$ is fulfilled for clusters of any size if we demand $C_l < 1/n$; this is really very mild because the clusters are chosen in practice according to (28) with $C = 1$, i.e. $C_l = \epsilon\kappa(R') \max(\kappa(X), \kappa(Y))$, which is smaller than $1/n$ for moderate values of $\max(\kappa(X), \kappa(Y))$. This has led us to set in the numerical experiments

$$C_l = \min\{\epsilon\kappa(R') \max(\kappa(X), \kappa(Y)), 1/n\}. \quad (56)$$

With this choice the factor $1/(1 - (d_1 - 1)C_l)$ is always less than n , but it is just a little greater than 1 when $C_l \approx \epsilon$. The presence of the ratio $C_l/rg(\Sigma_1)$ may look odd because we are

⁶Notice that notation similar to the definition 5.2 has already been used in the Introduction (see (4) and (8)).

bounding precisely the quotient $\min\{rg(\Lambda_1^+), rg(\Lambda_1^-)\}/rg(\Sigma_1)$, however notice that Definition 5.1 and Lemma 5.3 imply

$$C_l < rg(\Sigma_1) \quad \text{and} \quad C_l < relgap(\Sigma_1). \quad (57)$$

The ratio $C_l/rg(\Sigma_1)$ is kept in the bound because $C_l \ll rg(\Sigma_1)$ may often happen. It is convenient to bear in mind that these remarks also hold for the bounds appearing in the next theorems of this section.

Now we consider a signed cluster whose relative closest cluster has at least one singular value corresponding to an eigenvalue with the same sign. In this situation, next theorem shows that the singular value relative gap is equivalent to the eigenvalue relative gap up to a moderate constant.

Theorem 5.6 *Let Σ_1 be a cluster of singular values and Σ_2 its relative closest cluster having d_2 elements, both of tolerance C_l . Let all the elements of Λ_1 have the same sign and at least one element of Λ_2 have the same sign as those of Λ_1 . If $(d_2 - 1)C_l < 1$, then*

$$rg(\Lambda_1) \leq \left(1 + \frac{2}{1 - (d_2 - 1)C_l} \frac{(d_2 - 1)C_l}{relgap(\Sigma_1)}\right) rg(\Sigma_1).$$

Theorems 5.5 and 5.6 guarantee that, in order to obtain (54) for all the singular value clusters, we only need to deal with signed clusters whose relative closest cluster is oppositely signed. This will be the setting for the rest of the section. The following theorem proves that under mild conditions joining clusters of this kind leads to (54).

Theorem 5.7 *Let Σ_1 be a cluster of d_1 elements and Σ_2 its relative closest cluster, having d_2 elements, both of tolerance C_l . Suppose that all the elements of Λ_1 have the same sign, and all the elements of Λ_2 have the opposite sign. Moreover assume that $(d - 1)C_l < 1$, where $d = \max\{d_1, d_2\}$. If $rg(\Sigma_1) < t < 1$ and*

$$rg(\Sigma_1 \cup \Sigma_2) > \min\{rg(\Sigma_1), rg(\Sigma_2)\} \quad (58)$$

then

$$\min\{rg(\Lambda_1), rg(\Lambda_2)\} \leq \frac{1}{1 - t} \left(1 + \frac{1}{1 - (d - 1)C_l} + \frac{1}{1 - (d - 1)C_l} \frac{(d - 1)C_l}{rg(\Sigma_1 \cup \Sigma_2)}\right) rg(\Sigma_1 \cup \Sigma_2).$$

The assumption $rg(\Sigma_1) < t < 1$ means that only singular value clusters whose relative gap are small enough need to be joined to other cluster in order to obtain (54). In practice we have set $t = relgap(\Lambda_1)/2$. Therefore, if $rg(\Sigma_1) \geq t$, the bound in Theorem 4.7 leads trivially to (54). The assumption (58), $rg(\Sigma_1 \cup \Sigma_2) > \min\{rg(\Sigma_1), rg(\Sigma_2)\}$, is imposed to guarantee that by joining clusters Σ_1 and Σ_2 when computing bases of invariant subspaces some improvement is achieved in the bound in Theorem 4.7. In this regard one may wonder what happens with $\max\{rg(\Sigma_1), rg(\Sigma_2)\}$, i.e. how much can the bound (51) worsen for the cluster with the maximum relative gap when Σ_1 and Σ_2 are joined? Next Lemma shows that no significant worsening may happen.

Lemma 5.8 *If both (58) and $rg\{\Sigma_1\} < t < 1$ are fulfilled, then*

$$\max\{rg(\Sigma_1), rg(\Sigma_2)\} < \frac{rg(\Sigma_1 \cup \Sigma_2)}{1 - t}$$

Notice that the difference between the maximum and the minimum of $\{rg(\Sigma_1), rg(\Sigma_2)\}$ in this case is again a consequence of the lack of symmetry of the relative error.

In order to obtain (54) for all the clusters we have to prove that if Σ_1 and its relative closest cluster Σ_2 , defined as in Theorem 5.7, do not fulfill (58) they will not be joined because Σ_1 has a singular value relative gap not worse, up to a moderate constant, than either its eigenvalue relative gap or the eigenvalue relative gap of Σ_2 . This is the goal of the next theorem.

Theorem 5.9 *Let Σ_1 be a cluster of d_1 elements and Σ_2 its relative closest cluster, having d_2 elements, both of tolerance C_l . Suppose that all the elements of Λ_1 have the same sign, and all the elements of Λ_2 have the opposite sign. Moreover assume that $(d-1)C_l < 1$, where $d = \max\{d_1, d_2\}$. If $rg(\Sigma_1) < t < 1$ and*

$$rg(\Sigma_1 \cup \Sigma_2) = \min\{rg(\Sigma_1), rg(\Sigma_2)\}, \quad (59)$$

then

$$\min\{rg(\Lambda_1), rg(\Lambda_2)\} \leq \frac{1}{1-t} \left(1 + \frac{1}{1-(d-1)C_l} + \frac{1}{1-(d-1)C_l} \frac{(d-1)C_l}{rg(\Sigma_1)} \right) rg(\Sigma_1).$$

Observe that hypothesis (59) is simply the negation of (58) because always $rg(\Sigma_1 \cup \Sigma_2) \geq \min\{rg(\Sigma_1), rg(\Sigma_2)\}$.

Although similar, the bounds appearing in Theorems 5.7 and 5.9 are different in the following sense. While in Theorem 5.7 $\min\{rg(\Lambda_1), rg(\Lambda_2)\} \approx rg(\Sigma_1 \cup \Sigma_2)$ holds always, in Theorem 5.9 $\min\{rg(\Lambda_1), rg(\Lambda_2)\} \ll rg(\Sigma_1)$ might happen. Thus the error bounds obtained replacing in (51) $rg(\Sigma_1)$ by $\min\{rg(\Lambda_1), rg(\Lambda_2)\}$ may be pessimistic in the conditions of Theorem 5.9.

Our last result shows that in order to obtain (54) unions of more than two clusters are not necessary. In the following Theorem three clusters are considered. Two of them satisfy the assumptions of Theorem 5.7, and the third cluster may be a candidate to join to the others. In this situation it will be proved that the relative singular value gap for the third cluster is equivalent, up to a moderate constant, to its eigenvalue relative gap.

Theorem 5.10 *Let Σ_1 and Σ_2 be clusters satisfying the hypotheses of Theorem 5.7. Let Σ_3 be another cluster, of tolerance C_l , with all the elements of Λ_3 of the same sign and $rg(\Sigma_3) < t_3 < 1$. If Σ_1 (resp. Σ_2) is the relative closest cluster to Σ_3 , and all the elements of Λ_3 have sign opposite to those of Λ_1 (resp. Λ_2), then*

$$rg(\Lambda_3) \leq \left(1 + \frac{1}{(1-t)(1-t_3)} \frac{1}{1-(d-1)C_l} + \frac{1+t_3}{1-(d-1)C_l} \frac{(d-1)C_l}{rg(\Sigma_3)} \right) rg(\Sigma_3).$$

As announced after Definition 5.2, all the bounds appearing in this section remain true if every rg is replaced by the corresponding $relgap$. This is easily understood as follows: the left hand sides of the inequalities decrease if the rg 's are replaced by the $relgap$'s, and the new left hand sides are smaller than or equal to 1. The factors which multiply the rg 's appearing in the right hand sides are all of them greater than or equal to 1, and increase when quotients of the kind C_l/rg are replaced by $C_l/relgap$. Thus the left hand sides are bounded simultaneously by 1 and by some factor greater than or equal to 1 times the corresponding rg . Then they are bounded by the factor times the $relgap$. Also notice that for testing the assumptions in the results in this section, it is equivalent to use rg 's or $relgap$'s. First, it is trivial to see that $rg(\Sigma_1) < t < 1$ if and only if $relgap(\Sigma_1) < t < 1$. Second, in testing the condition (58), the following elementary Lemma holds:

Lemma 5.11 *Let*

$$\Sigma_1 = \{\sigma_{i+1}, \sigma_{i+2}, \dots, \sigma_{i+d_1}\} \quad , \quad \Sigma_2 = \{\sigma_{i+d_1+1}, \sigma_{i+d_1+2}, \dots, \sigma_{i+d_1+d_2}\}$$

be any pair of consecutive clusters of nonzero singular values of tolerance C_l . Then

1. $rg(\Sigma_1 \cup \Sigma_2) = \min\{rg(\Sigma_1), rg(\Sigma_2)\}$ *if and only if*
 $relgap(\Sigma_1 \cup \Sigma_2) = \min\{relgap(\Sigma_1), relgap(\Sigma_2)\}$.
2. $rg(\Sigma_1 \cup \Sigma_2) > \min\{rg(\Sigma_1), rg(\Sigma_2)\}$ *if and only if*
 $relgap(\Sigma_1 \cup \Sigma_2) > \min\{relgap(\Sigma_1), relgap(\Sigma_2)\}$.

The key fact to prove this simple Lemma is that $rg(\Sigma_1) \leq (\sigma_{i+d_1} - \sigma_{i+d_1+1})/\sigma_{i+d_1} < 1$, thus the 1 appearing in the *relgap*'s does not play any role. Taking into account that always $rg(\Sigma_1 \cup \Sigma_2) \geq \min\{rg(\Sigma_1), rg(\Sigma_2)\}$ and $relgap(\Sigma_1 \cup \Sigma_2) \geq \min\{relgap(\Sigma_1), relgap(\Sigma_2)\}$, either statement in previous lemma is just the negation of the other one.

The final consequence of this section is that in order to get (54) only clusters fulfilling the hypotheses of Theorem 5.7 must be joined. Once a pair of clusters of this kind are joined, they can be disregard in any other union processes as shown by Theorem 5.10. Otherwise, the rest of the results prove that union of clusters of different kind is not needed. In the next section the task of developing a routine that selects and joins clusters according this criterion will be undertaken.

5.2 Choosing a new set of clusters

Now we will present a routine for **step 3** of Algorithm 3. Given a set of clusters as input, selected according to (28), a new set of cluster will come out according to the logic of the theorems in subsection 5.1, i.e. clusters will be joined only if the hypotheses of Theorem 5.7 are satisfied.

Algorithm 3.1

Input: Eigenvalues: Λ ; Clusters: $\{\Sigma_i\}_{i=1}^k$; *tolgap*: parameter smaller than 1;.
Output: New set of clusters: $\{\Sigma_i\}_{i=1}^q$ with $q \leq k$.

Notation:

Λ_i denotes the set of eigenvalues whose absolute values are the elements of Σ_i .

1. $q = k$
 2. for $i=1:k$;
 - if $(\lambda_j > 0 \quad \forall \lambda_j \in \Sigma_i)$ then
 - $sign(\Sigma_i) = +1$
 - elseif $(\lambda_j < 0 \quad \forall \lambda_j \in \Sigma_i)$
 - $sign(\Sigma_i) = -1$
 - else
 - $sign(\Sigma_i) = 0$
 - $\frac{relgap(\Sigma_i)}{relgap(\Lambda_i)} = 2$
- endif

```

endfor
3.  $qrgmin = \min_{1 \leq i \leq q} \frac{relgap(\Sigma_i)}{relgap(\Lambda_i)} \equiv \frac{relgap(\Sigma_{i_c})}{relgap(\Lambda_{i_c})}$ 
4. while  $qrgmin < tolgap$ 
    determine the relative closest7 cluster  $\Sigma_{i_c+1}$  to  $\Sigma_{i_c}$  according to Def. 5.4
    if  $(sign(\Sigma_{i_c}) * sign(\Sigma_{i_c+1}) = -1)$  .and.
         $(relgap(\Sigma_{i_c} \cup \Sigma_{i_c+1}) > \min\{relgap(\Sigma_{i_c}), relgap(\Sigma_{i_c+1})\})$  then
             $q = q - 1$ 
             $relgap(\Sigma_{i_c}) = relgap(\Sigma_{i_c} \cup \Sigma_{i_c+1})$ 
             $sign(\Sigma_{i_c}) = 0$ 
             $\Sigma_{i_c} = \Sigma_{i_c} \cup \Sigma_{i_c+1}$ 
            for  $k = i_c + 1 : q$ 
                 $\Sigma_k = \Sigma_{k+1}$ 
                 $relgap(\Sigma_k) = relgap(\Sigma_{k+1})$ 
                 $sign(\Sigma_k) = sign(\Sigma_{k+1})$ 
            endfor
        endif
         $\frac{relgap(\Sigma_{i_c})}{relgap(\Lambda_{i_c})} = 2$ 
         $qrgmin = \min_{1 \leq i \leq q} \frac{relgap(\Sigma_i)}{relgap(\Lambda_i)} \equiv \frac{relgap(\Sigma_{i_c})}{relgap(\Lambda_{i_c})}$ 
5. endwhile

```

In practice we have set $tolgap = 1/2$, but other values are admissible. This choice leads to values $t = (relgap(\hat{\Lambda}_i)/2) \leq (1/2)$ for the parameters t appearing in Theorems 5.7, 5.9 and 5.10.

For the new set of clusters selected by Algorithm 3.1, the error in the corresponding bases of invariant subspaces computed by Algorithm 2.2 is given by Theorem 4.7 using the new singular value relative gaps, and these are the sharpest bounds we have for Algorithm 3. Nevertheless, in the next theorem we will use the theorems in the previous subsection to give an upper bound for the inverse of the new singular value relative gaps in (51) in terms of inverses of the eigenvalue relative gaps. Therefore this theorem ⁸ gives a precise statement of (54).

Theorem 5.12 *Let A be a $n \times n$ real symmetric matrix of rank r for which it is possible to compute a RRD fulfilling (9). Let $\hat{\Sigma}$ be the singular values of A computed using steps 1-2 of Algorithm 1. Let $\hat{\Sigma}_i$, $i = 1, \dots, q$, be the set of clusters of nonzero computed singular values of A selected by step 3 of Algorithm 3, $\hat{\Lambda}_i = \hat{\Lambda}_i^+ \cup \hat{\Lambda}_i^-$, $i = 1, \dots, q$, the corresponding set of clusters of eigenvalues and $\hat{Q}_i = [\hat{Q}_i^+ \hat{Q}_i^-]$, $i = 1, \dots, q$, the matrices computed by step 4 of Algorithm 3. Let Σ_i (resp. Λ_i), $i = 1, \dots, q$, be the corresponding clusters of exact singular values (resp. eigenvalues).*

⁷The same can be done if Σ_{i_c-1} is the relative closest cluster to Σ_{i_c} .

⁸In Algorithm 3.1 some fractions $relgap(\Sigma_i)/relgap(\Lambda_i)$ have been set equal to 2 to indicate either that the cluster is unsigned, or that it has already been analyzed. However, both the numerator and denominator keep their values (which would make the actual value of the fraction smaller than or equal to one). Notice that Theorem 5.12 is stated in terms of the actual values of $relgap(\hat{\Sigma}_i)$ and $relgap(\hat{\Lambda}_i)$.

1. If neither $\widehat{\Lambda}_i^+$ nor $\widehat{\Lambda}_i^-$ are empty, then there exist matrices Q_i^+ and Q_i^- , whose columns form orthonormal bases of the invariant subspaces of A corresponding, respectively, to the positive and negative eigenvalues of Λ_i , such that

$$\|\widehat{Q}_i^+ - Q_i^+\|_F \leq \frac{O(\epsilon\kappa(R') \max(\kappa(X), \kappa(Y)))}{\min\{\text{relgap}(\widehat{\Lambda}_i^+), \text{relgap}(\widehat{\Lambda}_i^-)\}}, \quad (60)$$

with a similar bound for $\|\widehat{Q}_i^- - Q_i^-\|_F$.

2. If all the elements of $\widehat{\Lambda}_i$ have the same sign and $\text{relgap}(\widehat{\Sigma}_i) \geq \mathbf{tolgap} * \text{relgap}(\widehat{\Lambda}_i)$, then there exist a matrix Q_i , whose columns form an orthonormal basis of the invariant subspace of A corresponding to the eigenvalues in Λ_i , such that

$$\|\widehat{Q}_i - Q_i\|_F \leq \frac{O(\epsilon\kappa(R') \max(\kappa(X), \kappa(Y)))}{\text{relgap}(\widehat{\Lambda}_i)}. \quad (61)$$

3. If all the elements of $\widehat{\Lambda}_i$ have the same sign, $\text{relgap}(\widehat{\Sigma}_i) < \mathbf{tolgap} * \text{relgap}(\widehat{\Lambda}_i)$, and the relative closest cluster $\widehat{\Sigma}_{cl(i)}$ to $\widehat{\Sigma}_i$ has all the corresponding eigenvalues with the opposite sign, then there exist a matrix Q_i , whose columns form an orthonormal basis of the invariant subspace of A corresponding to the eigenvalues in Λ_i , such that

$$\|\widehat{Q}_i - Q_i\|_F \leq \frac{O(\epsilon\kappa(R') \max(\kappa(X), \kappa(Y)))}{\min\{\text{relgap}(\widehat{\Lambda}_i), \text{relgap}(\widehat{\Lambda}_{cl(i)})\}}. \quad (62)$$

4. If all the elements of $\widehat{\Lambda}_i$ have the same sign, $\text{relgap}(\widehat{\Sigma}_i) < \mathbf{tolgap} * \text{relgap}(\widehat{\Lambda}_i)$, and the relative closest cluster to $\widehat{\Sigma}_i$ does not have all the corresponding eigenvalues with the opposite sign, then there exist a matrix Q_i , whose columns form an orthonormal basis of the invariant subspace of A corresponding to the eigenvalues in Λ_i , such that

$$\|\widehat{Q}_i - Q_i\|_F \leq \frac{O(\epsilon\kappa(R') \max(\kappa(X), \kappa(Y)))}{\text{relgap}(\widehat{\Lambda}_i)}. \quad (63)$$

Furthermore, let $\widehat{Q} = [\widehat{Q}_1^+ \ \widehat{Q}_1^- \ \dots \ \widehat{Q}_q^+ \ \widehat{Q}_q^-]$ be the $n \times r$ matrix whose columns are the bases of all considered invariant subspaces of A computed using **step 4** of Algorithm 3. Then there exists an $n \times r$ matrix B with exact orthonormal columns such that

$$\|\widehat{Q} - B\|_F = O(\epsilon). \quad (64)$$

Proof: The proof follows from Theorem 4.7 applied to the output clusters of Algorithm 3.1 (**step 3** of Algorithm 3) and the theorems on gaps in Section 5.1 with $C_l = \epsilon\kappa(R') \max(\kappa(X), \kappa(Y))$. As remarked following Theorem 5.10, *relgap*'s instead of *rg*'s can be used in the theorems of Section 5.1.

We begin by replacing $\text{relgap}(\Sigma_i, \widehat{\Sigma}_i)$ by $\text{relgap}(\widehat{\Sigma}_i, \Sigma_i)$ in the bound (51). This does not change significantly the bound due to (23). Moreover, we assume $\text{relgap}(\widehat{\Sigma}_i, \Sigma_i) \approx \text{relgap}(\widehat{\Sigma}_i, \widehat{\Sigma}_i)$ as was done right after (37). This is a fair assumption whenever steps 1-2 of Algorithm 1 compute singular values with high relative accuracy. Thus (55) allows to apply (51), with $\text{relgap}(\Sigma_i, \widehat{\Sigma}_i)$ replaced by $\text{relgap}(\widehat{\Sigma}_i)$, to the clusters selected by Algorithm 3.1.

Consider a cluster $\widehat{\Sigma}_{i_c}$ of singular values corresponding to the quantity *qrgmin* in Algorithm 3.1. This cluster is joined to its relative closest cluster if and only if the following three conditions are simultaneously fulfilled:

$$(c1) \frac{relgap(\widehat{\Sigma}_{i_c})}{relgap(\widehat{\Lambda}_{i_c})} < \mathbf{tolgap} < 1.$$

$$(c2) \mathit{sign}(\widehat{\Sigma}_{cl(i_c)}) * \mathit{sign}(\widehat{\Sigma}_{i_c}) = -1, \text{ where } \widehat{\Sigma}_{cl(i_c)} \text{ is the closest cluster to } \widehat{\Sigma}_{i_c}.$$

$$(c3) relgap(\widehat{\Sigma}_{i_c} \cup \widehat{\Sigma}_{cl(i_c)}) > \min\{relgap(\widehat{\Sigma}_{i_c}), relgap(\widehat{\Sigma}_{cl(i_c)})\}.$$

If all three conditions (c1), (c2) and (c3) are fulfilled, Algorithm 3.1 joins $\widehat{\Sigma}_{i_c}$ and $\widehat{\Sigma}_{cl(i_c)}$ in a new output cluster $\widehat{\Sigma}_{i_c} \cup \widehat{\Sigma}_{cl(i_c)}$. In this case Theorem 5.7 applies with $t = \mathbf{tolgap} * relgap(\widehat{\Lambda}_{i_c})$. This together with (51) yields (60) for the eigenvectors corresponding to the new output cluster.

Now, suppose that at least one of the three conditions is not satisfied. Suppose first that (c1) is satisfied, which implies $\mathit{sign}(\widehat{\Sigma}_{i_c}) \neq 0$, otherwise $qrgmin = 2$. If (c2) is not verified and the closest cluster to $\widehat{\Sigma}_{i_c}$ is an input cluster, Theorem 5.6 can be applied to the bound (51) to obtain (63); on the other hand if (c2) is not verified and the closest cluster is a new output cluster, (63) is achieved by using Theorems 5.6 or 5.10. If (c2) is verified and (c3) is also verified we are in the previously studied case of joining clusters. If (c2) is verified and (c3) is not verified, Theorem 5.9 can be applied to (51) to yield (62).

Suppose from now on that (c1) is not satisfied. Then, Algorithm 3.1 stops and all the clusters existing at that moment verify

$$\frac{relgap(\widehat{\Sigma}_i)}{relgap(\widehat{\Lambda}_i)} \geq \mathbf{tolgap}, \quad i = 1, \dots, q.$$

If $\mathit{sign}(\widehat{\Sigma}_i) = 0$, this is either because $\mathit{sign}(\widehat{\Sigma}_i) = 0$ on input, or because $\widehat{\Sigma}_i$ is a new output cluster, union of two input clusters. Anyway, Theorems 5.5 or 5.7 lead to (60) by using (51). If $\mathit{sign}(\widehat{\Sigma}_i) \neq 0$ and $relgap(\widehat{\Sigma}_i)/relgap(\widehat{\Lambda}_i) = 2$, then $\widehat{\Sigma}_i$ has been already analyzed inside the `while` loop and, according to the previous paragraph, either (62) or (63) are satisfied. If $\mathit{sign}(\widehat{\Sigma}_i) \neq 0$ but $\mathbf{tolgap} \leq relgap(\widehat{\Sigma}_i)/relgap(\widehat{\Lambda}_i) \leq 1$, then (51) implies (61), at the cost of an additional factor $1/\mathbf{tolgap}$. With this all the possible cases on the decision tree for the conditions (c1), (c2) and (c3) have been studied. The proof of (64) is as in Theorem 4.7. ■

We finish this section with two important remarks:

Remark 1: The *eigenvalue* clusters treated in the last theorem are exactly the same as the ones corresponding to the singular value clusters chosen according to (28). Therefore the bases treated in Theorem 5.12 and in Theorem 4.7 correspond exactly to the same invariant subspaces. This is because Algorithm 3.1 only joins oppositely signed clusters and Algorithm 2.2 computes the bases separately.

Remark 2: The bounds in Theorem 5.12 have been obtained in two stages: first, applying Theorem 4.7 to the new set of clusters produces a bound depending on singular value relative gaps. Then, this bound is majorized by other ones, depending on certain eigenvalue relative gaps. This second stage never worsens significantly the first bound, except in *case 3* of the statement of Theorem 5.12. Thus, the bound (62) may be pessimistic, because $\min\{relgap(\widehat{\Lambda}_i), relgap(\widehat{\Lambda}_{cl(i)})\}$ might be much smaller than $relgap(\widehat{\Sigma}_i)$. However, recall that the sharpest bound for Algorithm 3 is $O(\epsilon\kappa(R') \max(\kappa(X), \kappa(Y)))/relgap(\widehat{\Sigma}_i)$.

6 Numerical experiments

In this section we present results of two types of numerical experiments: first, we test Algorithms 2 and 3, the two implementations of **step 3** in Algorithm 1. A second kind of experiments tests the whole Algorithm 1, including the computation of the rank-revealing decomposition in two different ways: either a symmetric RRD of the form $A = XDX^T$, or a nonsymmetric one of the form $A = XDY^T$. Besides, a first subsection describes some practical details of the implementation of the three steps of Algorithm 1.

As will be seen from the experiments in § 6.2, Algorithm 1 behaves as predicted by the error analysis in Sections 4 and 5, and compares well both in the sense of accuracy and of speed with the J-orthogonal algorithm.

6.1 Implementation of Algorithm 1 (SSVD)

1. The rank revealing decomposition of the matrix A in **step 1** of Algorithm 1 has been done in two ways:
 - Symmetric RRD, $A = XDX^T$, using a modification of the symmetric indefinite Bunch & Parlett (BP) decomposition [3]: more specifically, we have used an adapted version of the routine `SGJGT` in [23].
 - A non-symmetric RRD, $A = XDY^T$, by means of a LU factorization with complete pivoting (GECV). We have used a modification of the LAPACK procedure `SGETF2`.
2. The singular value decomposition in **step 2** of Algorithm SSVD has been done using Algorithm 3.1 of [7]. Only LAPACK and BLAS routines have been used, as in [7], except for the one-sided Jacobi code in which we have used a routine developed by Z. Drmač according to the ideas in [11]. The implementation of the procedure (called `SGEPSV` in [7] in single precision) has followed the steps:

Algorithm 4 (SGEPSV) (Algorithm 3.1 in [7])

Input: $X, D, Y : A = XDY^T$. Output: $U, \Sigma, V : A = U\Sigma V^T$.

- (a) QR factorization with column pivoting of XD , $XD P = QR$; $A = QR P^T Y^T$
LAPACK Routine: `SGEQPF`
- (b) Multiply to get $W = R(Y P)^T$; $A = QW$
BLAS Routine: `STRMM`
- (c) SVD of W with one-sided Jacobi; $W = \bar{U}\Sigma V^T$; $A = Q\bar{U}\Sigma V^T$
Routine: `S_SGESVDJ` developed by Z. Drmač [11]
- (d) Multiply $U = Q\bar{U}$; $A = U\Sigma V^T$
LAPACK Routine: `SORMQR`

Two versions of this algorithm have been used, depending on whether right-Jacobi (right multiplication on W by Jacobi plane rotations) or left-Jacobi (right multiplication on W^T by Jacobi plane rotations) is employed in the one-sided Jacobi step. The left-Jacobi version has the advantage of speeding up the convergence. Although the error bounds for this version are weaker than for the other one (see Appendix A), no significant difference in accuracy has ever been observed in practice. Our experiments in § 6.2 confirm this.

In any case the routine that has been used computes one of the singular vector matrices by a product of Jacobi plane rotations. There exist much faster, equally accurate, versions of one-sided Jacobi algorithms which do not accumulate rotations [13], and could also be used. Nevertheless, with the present implementation the timing statistics of Algorithm SSVD is comparable to the J-orthogonal algorithm (see the timing data in the last paragraph of *Experiment 2* in subsection 6.2).

3. Algorithm 2 in **step 3** of Algorithm 1 has been implemented as described in subsection §3.3. The alternative Algorithm 3 in **step 3** of Algorithm SSVD has been implemented as described in section §5.

Some more specific details are the following:

- (a) Recall that **steps 1,2** are the same in both algorithms 2 and 3, and therefore the eigenvalues calculated by both algorithms are the same.
- (b) The choice of clusters in **step 1** of Algorithms 2 and 3 has been done using (28) by taking $C = 1$ and using the $O(n^2)$ estimator routine STRCON in LAPACK to estimate $\kappa(R')$, or $\kappa(X)$, $\kappa(Y)$ when starting from a non-factorized matrix. When generating matrices in RRD form $A = XDX^T$, some matrices X producing values of $\epsilon\kappa(R')\kappa(X)$ larger than one have appeared. This means that Algorithm 4 guarantees no significant digits of precision in the computation of the singular values. Moreover, using (28) produces in this case that all singular values are contained in just one cluster. Our discussion after Theorem 5.5 has led us to establish in practice the criterion to include two contiguous singular values σ_j, σ_{j+1} in the same cluster whenever

$$\frac{|\sigma_j - \sigma_{j+1}|}{\sigma_j} \leq \min\{\epsilon\kappa(R') \max\{\kappa(X), \kappa(Y)\}, 1/n\}. \quad (65)$$

- (c) The product $\Delta_i = V_i^T U_i$ in **step 11** of Algorithm 2.2 has been done using the BLAS routine SGEMM.
- (d) The diagonalization of $\Delta_i = [W_i^+ W_i^-] J_i [W_i^+ W_i^-]^T$ (**step 12** of Algorithm 2.2) has been done using the LAPACK routine SSYEV applied only to the triangular upper half of the matrix, as assumed in Lemma 4.5. Finally the eigenvector matrices $Q_i^\pm = V_i W_i^\pm$ (**step 13** of Algorithm 2.2) are obtained using the BLAS multiplication routine SGEMM.
- (e) In all the experiments the value for the parameter `tolgap` appearing in Algorithm 3.1 has been set to `tolgap = 1/2`.

6.2 Numerical results

The experiments were done using an AMD K7 ATHLON processor with IEEE arithmetic, and the routines implemented with Fortran PowerStation 4.0TM, from Microsoft. All numerical experiments in this section have been done with nonsingular matrices, although as pointed out in Sections §3 and §4, Algorithm 1 SSVD can be also applied to rank-deficient matrices.

In the first experiment we start from matrices already in factorized RRD form $A = XDX^T$, directly generating the matrices X and D . This has allowed us to focus on the accuracy of **step 3** in Algorithm 1, since, given the RRD, the work by Demmel et al. in [7] allows to control the error in **step 2** of Algorithm 1.

In the second group of experiments, two different kinds of non-factorized test matrices have been generated: graded matrices, and matrices specifically designed in [23] to guarantee a good

performance of the J-orthogonal algorithm. The reason for choosing graded matrices is that it is known, under the conditions given in [7, Section 4], that an accurate RRD, in the sense of (9), can be computed using a *plain implementation* of Gaussian elimination with complete pivoting (GECF). For the rest of the classes of matrices treated in [7, pp. 26-27], special implementations of GECF are needed to get the desired accuracy, and it is unfair to compare with the J-orthogonal algorithm, since at present no special implementations of the symmetric indefinite factorization are known to guarantee the accuracy. The reason for choosing the matrices designed in [23] is to compare Algorithm 1 and the J-orthogonal algorithm on matrices where the accuracy of the latter is guaranteed.

To test Algorithm 1 we have used as reference the eigenvalues and eigenvectors computed by the routine `DSYEVJ`, developed by I. Slapničar, that implements the implicit one-sided J-orthogonal algorithm⁹ [23] in double precision ($\epsilon = \epsilon_D \approx 1.11 \times 10^{-16}$). From now on these eigenvalues and eigenvectors are denoted, respectively, simply by λ_i and q_i . These are compared with the eigenvalues and eigenvectors, $\lambda_i^{(S)}$ and $q_i^{(S)}$, computed in single precision ($\epsilon = \epsilon_s \approx 5.96 \times 10^{-8}$) by the following routines: `SSVD` (Algorithm 1, using Algorithm 2 `SYSSV` in **Step 3**), `SSVDR` (Algorithm 1, using Algorithm 3 `SYSSVR` in **Step 3**), `SSYEVJ` (J-orthogonal algorithm, denoted simply by J-O in tables and figures), and, only when we start from a full (not already in rank-revealing form) matrix A , `SJAC` (standard Jacobi algorithm with new the stopping criterion introduced in [8, pp. 1206] with `tol` = ϵ_s) and `SSYEV` (LAPACK driver routine that implements tridiagonalization followed by QR iteration). For these methods the following quantities have been measured for each test matrix:

1. The maximum relative error in the eigenvalues

$$e_\lambda^{(S)} = \max_i \left| \frac{\lambda_i - \lambda_i^{(S)}}{\lambda_i} \right|. \quad (66)$$

2. A control quantity for eigenvalues

$$\vartheta^{(S)} = \frac{e_\lambda^{(S)}}{\kappa \epsilon_s} \quad (67)$$

where $\kappa = \kappa(R') \max\{\kappa(X), \kappa(Y)\}$, as in (14). Observe that when referring to symmetric RRDs κ is just $\kappa(R')\kappa(X)$. According to the bound (38), the quantity $\vartheta^{(S)}$ is expected to be $O(1)$ for Algorithm `SSVD`. For the J-orthogonal algorithm the error $e_\lambda^{(S)}$ is essentially bounded by $O(\epsilon_s \kappa(XD_X))$, where XD_X is the best conditioned column diagonal scaling of matrix X [23]. However, we have checked that $\kappa(X) \approx \kappa(XD_X)$ in our tests. This is due to the fact that the matrices X appearing in our experiments do not have any special structure. Furthermore the extra factor $\kappa(R')$ in the denominator, that we have observed that is $O(n)$ in the numerical tests in this section (see also [7, Thm. 3.2]), renders $\vartheta^{(S)}$ inadequate to check how well the bounds for the J-orthogonal algorithm behave, though it is still valid to compare the accuracy of `SSVD` and the J-orthogonal algorithm. For the other two considered algorithms, Jacobi and QR, $\vartheta^{(S)}$ is just the maximum error in the eigenvalues normalized in the same way as for the `SSVD` and J-orthogonal algorithms. Similar remarks apply to the eigenvector computations.

⁹ `DSYEVJ` is a driver routine formed by two routines that implement the two steps of the J-orthogonal algorithm: subroutine `DGJGT` (symmetric indefinite decomposition with complete pivoting) and subroutine `DJGJF` (implicit J-orthogonal Jacobi method with the same stopping criterion as one-sided Jacobi). `DSYEVJ` has been used when starting with the full matrix A . When starting from a factorized matrix $A = XDX^T$ only the subroutine `DJGJF` has been used. Similar remarks apply to the single precision driver routine `SSYEVJ`.

3. Corresponding to each cluster of eigenvalues, the sine of the maximum of canonical angles between the subspaces spanned by the computed basis, Q_i , in double precision and the computed basis, $Q_i^{(S)}$, in single precision:

$$E_{\Lambda_i}^{(S)} = \|\sin \Theta(Q_i, Q_i^{(S)})\|_2, \quad (68)$$

In the case of clusters with one single element we have calculated just the euclidean norm of the difference between the computed eigenvectors in double, q_i , and in single precision, $q_i^{(S)}$,

$$e_{q_i}^{(S)} = \|q_i - q_i^{(S)}\|_2. \quad (69)$$

Actually, these quantities are always computed, even in the presence of clusters of dimension larger than one. We do this in order to check that clusters are only chosen whenever no accuracy can be guaranteed for individual computed eigenvectors.

4. The control quantities for bases of invariant subspaces are

$$\Xi_{\Sigma}^{(S)} = \max_i \frac{E_{\Lambda_i}^{(S)} \text{relgap}(\Sigma_i^{(S)})}{\kappa \epsilon_s}, \quad \Xi_{\Lambda}^{(S)} = \max_i \frac{E_{\Lambda_i}^{(S)} \text{relgap}(\Lambda_i^{(S)})}{\kappa \epsilon_s}, \quad (70)$$

and the corresponding ones for individual eigenvectors,

$$\xi_{\sigma}^{(S)} = \max_i \frac{\|q_i - q_i^{(S)}\|_2 \text{relgap}(\sigma_i^{(S)})}{\kappa \epsilon_s}, \quad \xi_{\lambda}^{(S)} = \max_i \frac{\|q_i - q_i^{(S)}\|_2 \text{relgap}(\lambda_i^{(S)})}{\kappa \epsilon_s}. \quad (71)$$

According to Theorem 4.7, $\Xi_{\Sigma}^{(S)}$ and $\xi_{\sigma}^{(S)}$ are expected to be $O(1)$ for Algorithms **SSVD** and **SSVDR**. Also $\Xi_{\Lambda}^{(S)}$ and $\xi_{\lambda}^{(S)}$ are expected to be $O(1)$ for the J-orthogonal algorithm, but not for Algorithms **SSVD** or **SSVDR**, because the accuracy of **SSVDR** is governed by Theorem 5.12. However, the quantities $\Xi_{\Lambda}^{(S)}$ and $\xi_{\lambda}^{(S)}$ will be computed by **SSVD** and **SSVDR** to check in practice how the **SSVDR** algorithm improves the accuracy of **SSVD** and how it compares with the J-orthogonal algorithm. Notice that the quantities $\text{relgap}(\Sigma_i^{(S)})$ correspond either to the set of cluster chosen according to (65) for Algorithm **SSVD**, or to the output clusters of Algorithm 3.1 for Algorithm **SSVDR**. The quantities $\text{relgap}(\Lambda_i^{(S)})$ are always the same, because the clusters for eigenvalues do not change (see remarks at the end of subsection 5.2). The *relgaps* in (71) are the ones defined in (4) and (8) for any of the algorithms.

For the sake of brevity, values of $\xi_{\sigma}^{(S)}$ or $\xi_{\lambda}^{(S)}$ are not shown for routines **SJAC** and **SSYEV**; we simply report that extremely large errors are obtained for these algorithms.

To do our experiments we have generated matrices in single precision in different ways. All the random matrices needed have been generated using the **LAPACK** routines **SLATM1**, for diagonal matrices, and **SLATMR**, for full matrices. When we have generated matrices with a fixed condition number \mathcal{K} , it has been done by producing diagonal matrices with elements of absolute values in the range from 1 to $1/\mathcal{K}$, and after that multiplying by random single precision orthogonal matrices. The distribution of the diagonal elements is controlled by the parameter **MODE** of the routine **SLATM1**: $|\text{MODE}| = 3$, geometrically distributed; $|\text{MODE}| = 4$, arithmetically distributed; $\text{MODE} = 5$, with logarithms uniformly distributed. If **MODE** is positive (resp. negative) the elements are set in decreasing (resp. increasing) order.

Experiment 1 This experiment is designed to test Algorithms 2 and 3. We have generated $n \times n$ matrices X and D (diagonal), factors of a matrix $A = XDX^T$, as done in [7]. Parameters have been chosen as follows: $\kappa(X) = 10^{[2:1:6]}$; $\kappa(D) = 10^{[2:2:16]}$; $MODE_X = 3, 4, 5$; $MODE_D = \pm 3, \pm 4, 5$. For each set of parameters we have run: 20 matrices for $n = 50, 100$ (total 12000 matrices for each n), 2 for $n = 250$ (total 1200 matrices), 2 for $n = 500$ (total 1200 matrices) and 1 for $n = 1000$, and only for 2 combinations of the $MODE_s$, (total 80 matrices).

Figure 1 shows the maximum, minimum and average (over all $MODE_s$, samples and $\kappa(D)$ s) of the quantity $\log_{10} e_\lambda^{(S)}$, roughly the number of correct digits in the computed eigenvalues, as a function of $\kappa(X)$ for $n = 100$ for Algorithm 1 (SSVD or SSVDR) and for the J-orthogonal algorithm. The line $\epsilon_s \kappa(X) \kappa(R')$ is plotted as a guide to the eye; the quantity $\kappa(R')$ in this line is really the average of $\kappa(R')$ over all the matrices with that value of $\kappa(X)$. The results confirm the theoretical error bounds for eigenvalues.

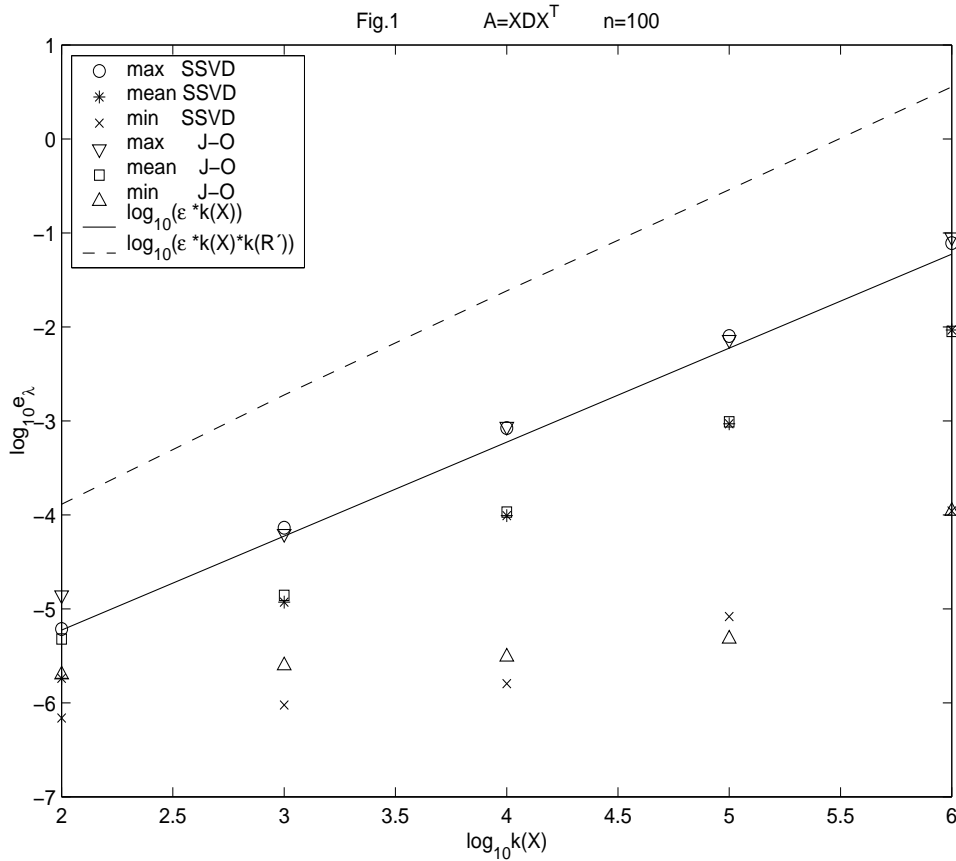


Figure 1: Experiment 1. Maximum relative error for eigenvalues: $\log_{10} e_\lambda^{(S)}$ against $\log_{10} \kappa(X)$.

Table 1 shows the statistical data corresponding to the quantity $\vartheta^{(S)}$. The aim is to check the bound (38) for Algorithm SSVD and compare its accuracy against the J-orthogonal algorithm. The most significant data in Table 1 appear under the columns named 'max' where the maximum values of each magnitude (the ones bounded by the error analysis) are shown. In particular, the fact that the quantities in the first row are smaller than 1 confirms that Algorithm 1 satisfies the bound (38). Besides, the third row itself is the control quantity ϑ calculated for the singular values computed in step 2 of Algorithm 1. The comparison of rows first and third shows that

n	50		100		250		500		1000	
Method	mean	max	mean	max	mean	max	mean	max	mean	max
$\vartheta^{(\text{SSVD})}$.030	.40	.022	.31	.015	.17	.013	.22	.013	.20
$\vartheta^{(\text{J-O})}$.041	.58	.037	.44	.039	.47	.044	.63	.050	.65
$\vartheta^{(\text{SVD})}$.030	.40	.022	.31	.015	.17	.013	.22	.012	.20

Table 1: Experiment 1. Statistical data for accuracy in eigenvalues: $\vartheta^{(S)}$.

Algorithm SSVD never misses a sign and always gives eigenvalues with the same precision as the singular values, except for five matrices of dimension 1000. These cases have $\kappa(X) = 10^6$ and $\epsilon_s \kappa(X) \kappa(R')$ greater than 100. Therefore *whenever $\epsilon_s \kappa(X) \kappa(R') < 1$ Algorithm SSVD has given the eigenvalues with the same precision as the singular values computed by Algorithm 3.1 in [7].* It can be seen, both from Figure 1 and Table 1, that Algorithm SSVD performs for eigenvalues as well (even a little better, specially for small values of $\kappa(X)$) as the J-orthogonal algorithm, with the maximum errors in Algorithm 1 adjusting very well to the predicted behavior $\epsilon \kappa(X) \kappa(R')$. It can be observed also that the data do not depend on n .

Moreover, for a significant portion of all the matrices (4144 matrices out of 12000 for $n = 50$; 6693 matrices out of 12000 for $n = 100$; 974 matrices out of 1200 for $n = 250$; 1105 matrices out of 1200 for $n = 500$; 79 matrices out of 80 for $n = 1000$), clusters of singular values of dimension greater than 1, according to criterion (65), have been found, with the maximum dimension of a cluster being 5. The average number of clusters has ranged from almost no clusters for $n = 50$, to around 40 clusters for $n = 1000$, with a typical dimension of 2. This shows that criterion (65) chooses clusters which determine perfectly in practice the signs of the eigenvalues. After applying Algorithm 3.1 all the considered matrices have clusters. The average number of clusters is in this case around $0.3n$ for all n .

In Table 2 we show the statistics for the number of orthogonal Jacobi sweeps for the Algorithm SSVD and the number of hyperbolic Jacobi sweeps for the J-orthogonal algorithm

n	50		100		250		500		1000	
Method	mean	max	mean	max	mean	max	mean	max	mean	max
$Sweeps^{(\text{SSVD})}$	5.5	10	6.3	12	7.4	12	8.4	14	9.3	15
$Sweeps^{(\text{J-O})}$	10.5	20	11.7	22	13.0	22	13.9	24	13.1	24

Table 2: Experiment 1. Statistical data for the number of sweeps.

These data correspond to the use of left-Jacobi in **step (c)** of Algorithm 4. If right-Jacobi is used the average number of sweeps for Algorithm SSVD is 13.8, with a maximum of 28 for $n = 100$, while the accuracy is the same. For these reasons, we have used in the rest of our experiments the left-handed version of the algorithm. It can be seen that the J-orthogonal algorithm uses more sweeps than the Algorithm SSVD: on average, from five more for $n = 50$ to almost four for $n = 1000$.

Now we focus on the analysis of data both for eigenvectors and for bases of invariant subspaces. Table 3 shows the quantities $\Xi_{\Sigma}^{(S)}$ and $\Xi_{\Lambda}^{(S)}$ defined in (70) for Algorithm 1, in both versions: SSVD, using Algorithm 2 SYSSV, and SSVDR, using Algorithm 3 SYSSVR. For the J-orthogonal algorithm we only show the quantity that governs its error: $\Xi_{\Lambda}^{(S)}$. When comparing the results of routines SSVD and SSVDR with the corresponding relative gaps of singular values

(rows 1 and 3), it can be seen that both methods behave as expected. When comparing the errors in the bases computed using the routine `SSVD` with the relative gap between eigenvalues the results can go rather badly (see row 2)¹⁰. When using `SSVDR` these results improve significantly (compare rows 4 and 2), showing that the method computes the bases for these test matrices with errors depending on the relative gap between eigenvalues, as the J-orthogonal algorithm does. It can be observed that the control quantities increase mildly with n for all the algorithms. Since this effect is not observed in the accuracy of the eigenvalues, this lead us to question if it is a real effect of the eigenvector bounds, or it is simply reflecting the fact that the quantities Ξ are computed from n -dimensional vectors.

n	50		100		250		500		1000	
Method	mean	max	mean	max	mean	max	mean	max	mean	max
$\Xi_{\Sigma}^{(SSVD)}$.032	.46	.051	1.2	.084	2.5	.12	4.5	.17	4.4
$\Xi_{\Lambda}^{(SSVD)}$.37	320	1.1	3300	2.4	500	6.5	1700	5.6	150
$\Xi_{\Sigma}^{(SSVDR)}$.034	.50	.056	1.2	.095	2.5	.13	4.5	.18	4.4
$\Xi_{\Lambda}^{(SSVDR)}$.041	.65	.075	4.6	.15	3.2	.23	6.0	.37	7.3
$\Xi_{\Lambda}^{(J-0)}$.044	.64	.076	1.5	.15	2.6	.21	5.7	.32	7.3

Table 3: Experiment 1. Statistical data for accuracy in bases of invariant subspaces: $\Xi_{\Sigma}^{(S)}$ and $\Xi_{\Lambda}^{(S)}$.

Table 4 shows the quantities $\xi_{\sigma}^{(S)}$ and $\xi_{\lambda}^{(S)}$ defined in (71). These are the quantities referring to the errors eigenvector by eigenvector. It can be seen that the accuracy of the eigenvectors is not spoiled by the clustering processes implicit in Algorithms `SSVD` and `SSVDR`. Similar comments to those made with respect to Table 3 apply here.

n	50		100		250		500		1000	
Method	mean	max	mean	max	mean	max	mean	max	mean	max
$\xi_{\sigma}^{(SSVD)}$.033	.74	.057	1.3	.092	2.5	.13	4.5	.19	4.4
$\xi_{\lambda}^{(SSVD)}$.37	320	1.1	3300	2.4	500	6.5	1700	5.6	150
$\xi_{\sigma}^{(SSVDR)}$.035	.90	.063	1.6	.10	2.5	.14	4.5	.20	4.4
$\xi_{\lambda}^{(SSVDR)}$.045	.90	.089	4.6	.17	3.2	.26	6.0	.42	7.3
$\xi_{\lambda}^{(J-0)}$.044	.64	.076	1.5	.15	2.6	.21	5.7	.32	7.3

Table 4: Experiment 1. Statistical data for accuracy in eigenvectors: $\xi_{\sigma}^{(S)}$ and $\xi_{\lambda}^{(S)}$.

To conclude, we show other quantities of numerical interest. The minimum singular value and eigenvalue relative gaps for clusters selected in Algorithm 2 have exceeded, respectively, 10^{-5} and 10^{-4} ; and after the clustering process in Algorithm 3 both relative gaps, for eigenvalues and singular values, have been bigger than 10^{-4} . The minimum relative gap for individual eigenvalues have been greater than 10^{-5} , and for singular values greater than 10^{-8} . The maximum values of $\kappa(R')$ have been: 190, for $n = 50$, 270, for $n = 100$, 600, for $n = 250$, 1300, for $n = 500$, 2200, for $n = 1000$, showing that it increases roughly as some constant times n .

¹⁰However, as can be deduced from the mean value of $\Xi_{\Lambda}^{(SSVD)}$, matrices for which `SSVD` computes eigenvectors with a large error with respect the relative gap between eigenvalues are quite infrequent.

Experiment 2 We have generated $n \times n$ graded matrices $A = DBD$ by multiplying random well-conditioned matrices, B , and random ill-conditioned diagonal matrices, D , to test the accuracy of the complete Algorithm 1 including the factorization in **step 1**. Not always can an accurate RRD fulfilling (9) be computed for graded matrices [7, §4]: the accuracy that can be guaranteed at best (and is frequently achieved in practice) is $O(\epsilon_s \kappa(B))$. Thus, high relative accuracy is expected when computing eigenvalues and eigenvectors for the matrices generated in this experiment. As mentioned in § 6.1, the initial rank revealing decomposition in Algorithms SSVD and SSVDR has been done in two ways: either using a modification of the symmetric indefinite Bunch & Parlett decomposition, or a non-symmetric LU factorization with complete pivoting. We have obtained similar results for both decompositions. Parameters have been chosen as follows: $\kappa(B) = 10^{[0:1:3]}$, $\kappa(D) = 10^{[2:2:10]}$, $MODE_B = 3, 4, 5$, $MODE_D = \pm 3, \pm 4, 5$. For each set of parameters we have run: 50 matrices for $n = 50, 100$ (total 15000 matrices for each n), 5 for $n = 250, 500$ (total 1500 matrices for each n), 1 for $n = 1000$, and only for 5 combinations of the $MODE$ s, (total 100 matrices). As announced, also Jacobi and QR have been applied on these test matrices.

The same quantities as in Experiment 1 are shown in Table 5 for eigenvalues, and in Table 6 for individual eigenvectors. The results for bases of invariant subspaces are almost the same as those in Table 6 and, therefore, are not shown. In these tables we show only the data corresponding to symmetric rank revealing decompositions obtained by the Bunch-Parlett method (abbreviated as BP in the table). The corresponding data for these tables using the unsymmetric rank revealing decomposition based on GECP are so similar that they are omitted. Nevertheless for other quantities (see Tables 7 and 8) we show the results for both decompositions (GECP is abbreviated as CP in the table).

Notice that the maximum values in Table 5 are greater than in Experiment 1, for both the SSVD and the J-orthogonal algorithm. This is due to the error in the initial factorization step, which is roughly bounded by $O(\epsilon_s \kappa(B))$. In any case, they behave much better than the classical methods, Jacobi and QR . An interesting remark is that the quantities $\vartheta^{(S)}$ decrease in Table 5 as n increases. This is because in this experiment (see Table 7) the condition number κ increases with the dimension n faster than the relative errors $e_\lambda^{(S)}$ in the eigenvalues. The control quantities for eigenvectors in Table 6 also decrease with n for the same reason. However, the maximum values of the control quantities for eigenvalues (Table 5) are much bigger than those of eigenvectors (Table 6). This is not explained by the error bounds.

n	50		100		250		500		1000	
Method	mean	max	mean	max	mean	max	mean	max	mean	max
$\vartheta^{(SSVD)}$	1.8	2600	.82	1100	.21	52	.22	140	.014	.24
$\vartheta^{(J-O)}$	1.5	1100	.80	1200	.21	64	.31	320	.019	.33
$\vartheta^{(JAC)}$	$3 \cdot 10^{15}$	$3 \cdot 10^{19}$	$1 \cdot 10^{14}$	$3 \cdot 10^{17}$	$1 \cdot 10^{13}$	$7 \cdot 10^{15}$	$7 \cdot 10^{12}$	$5 \cdot 10^{15}$	$2 \cdot 10^{11}$	$8 \cdot 10^{12}$
$\vartheta^{(QR)}$	$2 \cdot 10^{13}$	$2 \cdot 10^{17}$	$7 \cdot 10^{11}$	$5 \cdot 10^{15}$	$5 \cdot 10^{10}$	$4 \cdot 10^{13}$	$2 \cdot 10^{10}$	$1 \cdot 10^{13}$	$2 \cdot 10^3$	$4 \cdot 10^4$
$\vartheta^{(SVD)}$	1.8	2600	.82	1100	.21	52	.22	140	.014	.24

Table 5: Experiment 2. Statistical data for accuracy in eigenvalues: $\vartheta^{(S)}$.

As in Experiment 1, for a good number of the generated matrices (310 matrices out of 15000 for $n = 50$; 4821 matrices out of 15000 for $n = 100$; 1019 matrices out of 1500 for $n = 250$; 1454 matrices out of 1500 for $n = 500$; 100 matrices out of 100 for $n = 1000$), there are clusters of singular values of dimension greater than 1, according to criterion (65), with a maximal dimension of 5. The average number of clusters has ranged from almost no clusters for $n = 50$,

n	50		100		250		500		1000	
Method	mean	max	mean	max	mean	max	mean	max	mean	max
$\xi_\sigma^{(\text{SSVD})}$.47	11	.28	4.6	.17	1.1	.064	.55	.023	.16
$\xi_\lambda^{(\text{SSVD})}$	3.6	3300	2.8	1900	1.2	1600	.30	14	.067	.51
$\xi_\sigma^{(\text{SSVDR})}$.47	11	.31	5.2	.20	1.1	.076	1.3	.024	.16
$\xi_\lambda^{(\text{SSVDR})}$.56	12	.34	5.8	.25	2.4	.091	1.3	.030	.16
$\xi_\lambda^{(\text{J-0})}$.60	21	.37	4.3	.17	1.2	.090	.67	.039	.20

Table 6: Experiment 2. Statistical data for accuracy in eigenvectors: $\xi_\sigma^{(S)}$ and $\xi_\lambda^{(S)}$.

to around 60 clusters for $n = 1000$ with a typical dimension of 2. This shows again that criterion (65) determines perfectly in practice the signs of the eigenvalues, even when clusters are present. After applying Algorithm 3.1 all the considered matrices have clusters. The average number of clusters have been in this case around $0.3n$ for all n .

In addition, we show other quantities of numerical interest. The minimum singular value and eigenvalue relative gaps for clusters selected in Algorithm 2 are, respectively, 10^{-5} and $3.3 \cdot 10^{-4}$; and after the clustering process in Algorithm 3 both relative gaps, for eigenvalues and singular values, have reached the minimum $3.3 \cdot 10^{-4}$. The minimum relative gap for individual eigenvalues has been $4.1 \cdot 10^{-5}$, and for singular values greater than $9.1 \cdot 10^{-8}$. With respect to the condition numbers $\kappa(X)$, $\max\{\kappa(X), \kappa(Y)\}$ and $\kappa(R')$, they are shown in Table 7. The maximum values of $\epsilon\kappa(X)\kappa(R')$ are $8 \cdot 10^{-4}$, for $n = 50$, $4 \cdot 10^{-3}$, for $n = 100$, $5 \cdot 10^{-2}$, for $n = 250$, $3 \cdot 10^{-1}$, for $n = 500$, 1.8, for $n = 1000$, showing that it increases roughly as some constant times n .

n	50		100		250		500		1000	
Method	mean	max	mean	max	mean	max	mean	max	mean	max
$\kappa(R')$ (BP)	11	39	23	84	67	220	150	430	330	960
$\kappa(R')$ (CP)	11	37	24	80	71	201	160	450	360	860
$\kappa(X)$ (BP)	100	500	300	1300	1400	5000	4300	16000	14000	40000
$\max\{\kappa(X), \kappa(Y)\}$ (CP)	78	320	230	1000	1000	3200	2900	7900	5000	20000

Table 7: Experiment 2. Table for $\kappa(R')$ and $\max\{\kappa(X), \kappa(Y)\}$.

n	50		100		250		500		1000	
Method	mean	max	mean	max	mean	max	mean	max	mean	max
$Sweeps^{(\text{SSVD})\text{BP}}$	5.0	7	5.6	8	6.4	9	7.3	9	8.1	9
$Sweeps^{(\text{SSVD})\text{CP}}$	5.0	7	5.5	8	6.4	9	7.2	9	8.0	9
$Sweeps^{(\text{J-0})}$	6.3	8	7.1	10	8.5	11	9.6	12	11.0	13

Table 8: Experiment 2. Statistical data for the number of sweeps.

Table 8 shows that the J-orthogonal algorithm uses again more sweeps than the Algorithm SSVD: on average, from one more for $n = 50$ to three more for $n = 1000$. This is reflected in the run-time used by the different routines. Taking as a reference the time employed by the QR routine (SSYEV of LAPACK), we have the following average results for our experiments:

for $n = 100$, Algorithm **SSVDR** (with symmetric RRD factorization) employs 200% more time than QR, the J-orthogonal algorithm employs 250% more time and the Jacobi Algorithm **SJAC** employs 190% more time; for $n = 500$, Algorithm **SSVDR** (with symmetric RRD factorization) employs 380% more time, the J-orthogonal algorithm employs 350% more time and the Jacobi Algorithm **SJAC** employs 340% more time. This numbers can be explained as coming from two opposite effects: **SSVDR** uses less Jacobi sweeps but the number of clusters increases with the size of the matrix.

Experiment 3 We have also generated full matrices in another form to compare the accuracy of Algorithms 1 and J-orthogonal. We have used the matrix generator developed in [23], which is specifically designed to test the performance of the J-orthogonal algorithm on matrices for which the error bounds of this algorithm are controlled (see [23] for details).

The set of parameters have been chosen as follows: $n = 100$; **ASCAL** = [1 : 1 : 3]; **HSCAL** = [5 : 2 : 25]¹¹. For each set of parameters we have run 50 matrices: in total 1650 matrices.

The results confirm that Algorithm **SSVD** performs very well also for this type of matrices. The results for eigenvalues, eigenvectors and number of sweeps are summarized in Table 9. As in the other experiments, the results for individual eigenvectors, $\xi_{\sigma,\lambda}^{(s)}$, are similar to those for bases. For this set of matrices, no clusters of singular values with dimension greater than 1 were found in the sense of criterion (65).

n	ϑ		ξ_{σ}		ξ_{λ}		<i>Sweeps</i>	
Method	mean	max	mean	max	mean	max	mean	max
SSVDR	.27	2.2	2.1	14	2.9	21	4.6	6
J - 0	.47	2.8	—	—	3.1	20	5.5	8

Table 9: Experiment 3. Statistical data.

Experiment 4 The results for testing the accuracy of computed eigenvectors in previous experiments seem to show that the errors for the **SSVDR** and J-orthogonal algorithms are comparable (see rows 4 and 5 of Tables 4, 6 and Table 9 in Experiment 3), both depending on the relative gap between eigenvalues. However, it has not to be forgotten that the error bound for eigenvectors in the **SSVDR** algorithm is the expression (10) (or Theorem 5.12 for a more precise statement) and not (3). It is not difficult to think of situations in which Algorithm **SSVDR** can calculate single eigenvectors much worse than the J-orthogonal algorithm. Take for example the following 3×3 very well-conditioned matrix generated in single precision:

$$A = \begin{bmatrix} .1804019 & .9148742 & -.3611555 \\ .9148742 & -.2908984 & -.2799287 \\ -.3611555 & -.2799287 & -.8894936 \end{bmatrix}$$

with eigenvalues in double precision¹²: $\lambda_1 = 0.9999904633563307$, $\lambda_2 = -0.9999802814301686$, $\lambda_3 = -1.000000302456291$. The corresponding computed eigenvectors in single precision have the fol-

¹¹The routine **GENSYM** generates a non-singular symmetric matrix H of order n , with $\kappa(H) \approx 10^{\text{HSCAL}}$ and the measure $C(A, \hat{A}) \approx 10^{\text{ASCAL}}$ (see [23] for details).

¹²If the eigenvalues of matrix A are computed in MATLAB, the same numbers are not obtained. The reason is that the displayed matrix A is not exactly the matrix we used in our computations. A is just the rounding to seven decimal digits of its exact representation in single *binary* precision. However similar results are obtained.

lowing errors for the **SSVDR** algorithm

$$[\|q_i - q_i^{(\text{SSVDR})}\|_2]_{i=1,2,3} = [3.12, 5.25, 4.23] \times 10^{-3}$$

and

$$[\|q_i - q_i^{(\text{J-O})}\|_2]_{i=1,2,3} = [3.79 \times 10^{-5}, 1.43, 1.43] \times 10^{-3}$$

for the J-orthogonal algorithm. Notice that the J-orthogonal algorithm computes the eigenvector corresponding to the positive eigenvalue λ_1 with full machine precision, while with the **SSVDR** algorithm five significant decimal digits are lost. The reason of this is easily understood, because the eigenvalue relative gap for λ_1 is, approximately, 2, while the corresponding singular value relative gap is near 10^{-5} (in this case case relative or absolute gaps are equivalent). This cannot be improved by the clustering process done in Algorithm 3.1, because any of the two possible clusters of singular values containing one positive and one negative eigenvalue has a close singular value at a distance of order 10^{-5} , and the minimum of the eigenvalue relative gaps is also of order 10^{-5} .

However, notice that the **SSVDR** algorithm is able to compute all the eigenvectors with 3 correct decimal digits and that $\max_i e_{q_i}^{(\text{SSVDR})} / \max_i e_{q_i}^{(\text{J-O})} = 3.7$, of order 1 as predicted by bound (10), i.e. the J-orthogonal algorithm also computes some eigenvectors with 3 correct significant digits. One can easily modify the matrix A in such a way that **SSVDR** computes the eigenvector corresponding to the positive eigenvalue with an error of order 1, but then the J-orthogonal algorithm would lose all the correct digits in the eigenvectors corresponding to the negative eigenvalues.

Finally, notice that if all the eigenvalues of the matrix A are considered inside the same cluster, the **SSVDR** algorithm computes the eigenvector corresponding to λ_1 with full machine precision, according to the bound (51). However, the eigenvectors corresponding to the negative eigenvalues are computed with errors of order 1, although according to (51) they form a very accurate orthonormal basis of the invariant subspace associated with the negative eigenvalues.

Experiment 5 Our last experiment is designed to show how the **SSVD** algorithm, like the J-orthogonal one, is able to compute accurate bases of invariant subspaces, even when the gaps between eigenvalues are very small.

We generate a 10×10 matrix $A = QDQ^T$ by multiplying, in single precision, a single precision random orthogonal matrix Q by the diagonal matrix $D = \text{diag}[-1, 1, 1, 1, 1, 0.1, 0.1, 0.1, 0.1, 0.1]$. Due to roundoff errors, the absolute values of all the eigenvalues of A become different. But two clusters of singular values are found by Algorithm 2 (**SYSSV**) according to criterion (65), one around 1, of dimension 5, and another around 0.1, of the same dimension. The absolute gaps between the singular values inside each cluster exceed 10^{-7} . Thus the double precision routine **DSYEVJ** computes the eigenvectors with at least 8 correct decimal digits. The algorithms **SSVD** and J-orthogonal, in single precision, compute all the eigenvectors with errors of $O(1)$, except the eigenvector corresponding to the negative eigenvalue which is computed, in both cases, with an error near 10^{-7} . This error is predicted by bound (51) for the **SSVD** algorithm (see also the remarks after the proof of Theorem 4.7). The errors in the invariant subspaces can be estimated using $E_{\Lambda_i}^{(S)}$ in (68). These, for **SSVD** and J-orthogonal algorithms, have been of order 10^{-7} for the following invariant subspaces: the corresponding to the four positive eigenvalues close to 1, the corresponding to the five positive eigenvalues close to 0.1, and the corresponding to the negative eigenvalue. Moreover, the same error appear if we consider the invariant subspace corresponding to all the eigenvalues of absolute value around 1 (including the negative one). This shows in practice that, as studied in the error analysis leading to Theorem 4.7, once a cluster of singular

values is chosen, we obtain two bases, one for the invariant subspace corresponding to the positive eigenvalues in the cluster and another for the negative ones, with an error of the same order than the appearing in the bases of the singular subspaces corresponding to the whole cluster of singular values. In this experiment Algorithm 3.1 does not modify the set of clusters according to (28), since the clusters are not oppositely signed.

7 Conclusions and future work.

In this paper we have presented formal error analysis and numerical experiments of a new algorithm which computes eigenvalues and eigenvectors with high relative accuracy for the largest class of symmetric matrices known so far. In particular for all symmetric matrices belonging to the classes of general matrices studied in [7]. This high relative accuracy is achieved for a given symmetric matrix A whenever an accurate rank-revealing decomposition (RRD) of A can be computed.

The new algorithm is based on computing, in a first stage, a singular value decomposition of the symmetric matrix A . This is the reason of its wide applicability, because in this stage the symmetry of A is not used. Thus, we can compute non-symmetric RRDs of A and apply the theory developed in [7].

It is not known if accurate symmetric RRDs can be computed for all symmetric matrices in any of the classes described in [7]. The J-orthogonal algorithm [27, 23] computes eigenvalues and eigenvectors with high relative accuracy *only* if accurate enough symmetric RRDs are available. The authors are presently studying this interesting question.

A Appendix A: Backward error of the SVD algorithm

In this appendix we give a proof of Theorem 2.1, i.e. we show that the SVD algorithm employed in **step 2** of Algorithm SSVD produces a small backward multiplicative error when executed in finite precision arithmetic. More precisely, we analyze the following version for rectangular matrices of Algorithm 4 in § 6.1 (i.e., of Algorithm 3.1 in [7]). Recall that given a RRD XDY^T of a real m by n matrix G , $m \geq n$ of rank r , the inputs for Algorithm 4 are the three matrices $X \in \mathbb{R}^{m \times r}$, $D \in \mathbb{R}^{r \times r}$, $Y \in \mathbb{R}^{n \times r}$, with D diagonal and all three matrices of full rank r .

Algorithm 5 (Version of Algorithm 4 for rectangular matrices)

Input: rank-revealing decomposition XDY^T of $G \in \mathbb{R}^{m \times n}$.

Output: singular value decomposition $U\Sigma V^T$ of G .

1. compute a QR decomposition with column pivoting $XD = QRP$ of XD .
2. compute the product $W = RPY^T$ using conventional matrix multiplication.
3. compute a LQ decomposition $W = L_\omega Q_\omega^T$ of W .
4. compute an SVD $L_\omega = U_\omega \Sigma V_\omega^T$ of L_ω using right-handed Jacobi
5. compute the products $U = QU_\omega$ and $V = Q_\omega V_\omega$.

Notice that this implementation differs from the one presented both in §6.1 and in [7]: here the Jacobi step is split in two stages, steps 3 and 4. This is recommended in [7, §3.3] to save flops in the one-sided Jacobi computation, the most expensive one in the whole algorithm. This

only makes sense if the rank r is less than n . Otherwise, W is square and skipping step 3 above does not affect either the computational cost or the error bounds below. This is the case of the numerical experiments presented in Section 6.

The crucial ingredient to prove Theorem 2.1, which is missing in the analysis of [7], is that step 4 above, the one-sided Jacobi SVD algorithm on a square invertible matrix, produces a small multiplicative backward error. To be more specific, we will prove the following result, where the i -th column (resp. row) of \tilde{A} is denoted by $\tilde{A}(:, i)$ (resp. $\tilde{A}(i, :)$), and \tilde{A} is the last matrix in the sequence generated by the right-handed Jacobi process.

Theorem A.1 *Let $A \in \mathbb{R}^{n \times n}$ be an invertible matrix and let $\hat{U}\hat{\Sigma}\hat{V}^T$ be the SVD computed in finite arithmetic with machine precision ϵ by the right-handed Jacobi SVD algorithm on A with stopping criterion¹³*

$$\max_{i \neq j} \mathbf{fl} \left(\frac{|\tilde{A}(:, i)^T \tilde{A}(:, j)|}{\|\tilde{A}(:, i)\|_2 \|\tilde{A}(:, j)\|_2} \right) \leq n\epsilon, \quad \text{for } i \neq j. \quad (72)$$

Then, there exist matrices $U', V', E_L, E_R \in \mathbb{R}^{n \times n}$, such that U' and V' are orthogonal,

$$\begin{aligned} \|U' - \hat{U}\| &= O(\epsilon), & \|V' - \hat{V}\| &= O(\epsilon), \\ \|E_L\| &= O(\epsilon), & \|E_R\| &= O(\epsilon\kappa(A_N)), \end{aligned} \quad (73)$$

where A_N is the diagonal scaling of A with rows of unit euclidean norm and

$$(I + E_L)A(I + E_R) = U'\hat{\Sigma}V'^T. \quad (74)$$

Proof: It is known [12, Proposition 3.13] that, under the conditions above, the matrix \tilde{A} satisfying the stopping criterion (72) can be written as

$$\tilde{A} = (A + \delta A)V'$$

for an orthogonal matrix V' with $\|V' - \hat{V}\| = O(\epsilon)$ and δA such that

$$\|\delta A(i, :)\|_2 \leq \epsilon_J \|A(i, :)\|_2, \quad i = 1, \dots, n \quad (75)$$

for a certain $\epsilon_J = O(\epsilon)$ which depends on the sweeps required for convergence. Hence,

$$\tilde{A} = A(I + E_R)V' \quad (76)$$

for $E_R = A^{-1}\delta A$. If we now scale $A = D_N A_N$ with a diagonal matrix D_N so that A_N has rows of unit euclidean length, the bound (75) implies

$$\|E_R\|_F \leq \|A_N^{-1}\|_F \|D_N^{-1}\delta A\|_F \leq \sqrt{n}\epsilon_J \|A_N^{-1}\|_F,$$

where $\|\cdot\|_F$ stands for the Frobenius norm¹⁴. Finally, since $\|A_N\|_F = \sqrt{n}$, it follows that the Frobenius norm of E_R , and consequently its spectral norm, is bounded by $\epsilon_J \kappa_F(A_N) = O(\epsilon\kappa(A_N))$.

¹³A similar result holds with $n\epsilon$ replaced by any tolerance tol in criterion (72). In that case, $\|U' - \hat{U}\| \leq ntol + O(\epsilon)$ and $\|E_L\| \leq ntol + O(\epsilon)$. Notice, however, that if the tolerance is larger than $O(\epsilon)$ then the computed left singular vectors will fail, in general, to be orthogonal up to $O(\epsilon)$.

¹⁴One can also show that $\|E_R\| \leq \sqrt{n}\epsilon_J \|A_N^{-1}\|$ in the spectral norm.

On the other hand, recall that if we denote by $\tilde{\Sigma}$ the diagonal matrix whose i -th diagonal entry is the euclidean norm of the i -th column of \tilde{A} , then $\hat{\Sigma}$ and \hat{U} are computed as $\hat{\Sigma} = \mathbf{f1}(\tilde{\Sigma})$ and $\hat{U} = \mathbf{f1}(\tilde{A}\tilde{\Sigma}^{-1})$. Notice that each element \hat{u}_{ij} of \hat{U} can be written as $\hat{u}_{ij} = (\tilde{A}_{ij}/\tilde{\Sigma}_{jj})(1 + \epsilon_{ij})$ with $|\epsilon_{ij}| < \epsilon$. Let U be the matrix such that $\tilde{A} = U\hat{\Sigma}$. Then (76) implies that

$$U\hat{\Sigma}(V')^T = A(I + E_R)$$

with $\|U - \hat{U}\|_F \leq \epsilon\|U\|_F$. It only remains to show, using the stopping criterion, that there is an orthogonal matrix U' such that

$$U = (I + E_L)^{-1}U'$$

with $\|E_L\| = O(\epsilon)$ and $\|U' - \hat{U}\| = O(\epsilon)$.

It follows from condition (72) that each off-diagonal element of $U^T U$ is bounded in absolute value by $cn\epsilon + O(\epsilon^2)$, with c a small integer constant. The diagonal elements of $U^T U$, on the other hand, are $1 + \alpha$ with $|\alpha| \leq cn\epsilon + O(\epsilon^2)$. Thus, $\|U^T U - I\|_F \leq cn^2\epsilon + O(\epsilon^2)$. If $U = W_L(I + \delta\Sigma)W_R^T$ is an arbitrary SVD of U , then $\|\delta\Sigma\|_F \leq cn^2\epsilon + O(\epsilon^2)$. Denoting $U' = W_L W_R^T$, it follows that $U = (I + \delta U)U'$, where U' is orthogonal and $\|\delta U\|_F = \|\delta\Sigma\|_F$.

Defining $E_L = (I + \delta U)^{-1} - I$, we obtain that $\|E_L\|_F = \|\delta U\|_F + O(\|\delta U\|_F^2) \leq cn^2\epsilon + O(\epsilon^2)$.

Finally, $\|\hat{U} - U'\|_F \leq \|\hat{U} - U\|_F + \|U - U'\|_F$, but $\|U - U'\|_F = \|\delta U\|_F \leq cn^2\epsilon + O(\epsilon^2)$, and $\|\hat{U} - U\|_F \leq \epsilon\|U\|_F \leq \sqrt{n}\epsilon + O(\epsilon^2)$

■

We are now in the position to prove Theorem 2.1. Since we will cite results in [7, §3.2.1], we need to match our notation with that of [7]: the matrices Q, W, R (and R') appearing in the proof, which are the computed ones, are named in the proof *without* hats. The rest of the computed matrices are denoted, as elsewhere in this paper, with their hats on.

Proof of Theorem 2.1: It is shown in [7, p. 34] that, after step 2 of Algorithm 5, the matrix Q computed in step 1 and the matrix W computed in step 2 are such that

$$(I + E_1)G(I + F_1) = QW \tag{77}$$

for square matrices E_1, F_1 with

$$\|E_1\| = O(\epsilon\kappa(X)), \quad \|F_1\| = O(\epsilon\kappa(R')\kappa(Y)).$$

Although the computed Q has not exactly orthonormal columns, it is well known that there exists a matrix Q' with orthonormal columns such that

$$Q = Q' + E_q = (I + E_q(Q')^T)Q', \tag{78}$$

with $\|E_q\| = O(\epsilon)$. Thus, (77) becomes $(I + E_1')G(I + F_1) = Q'W$, with $\|E_1'\| = O(\epsilon\kappa(X))$.

The LQ factorization of W in Step 3 of Algorithm 5 is equivalent to computing a QR factorization of $W^T \in \mathbb{R}^{n \times r}$. The usual additive backward error analysis of the QR factorization, applied column-wise, ensures that the computed \hat{L}_ω satisfies

$$\hat{L}_\omega(Q'_\omega)^T = (W + E_\omega)$$

where $Q'_\omega \in \mathbb{R}^{n \times r}$ is a matrix with orthonormal columns satisfying $\|Q'_\omega - \hat{Q}_\omega\| = O(\epsilon)$ for the computed \hat{Q}_ω . The backward error E_ω satisfies the row-wise bound

$$\|E_\omega(i, :)\| = O(\epsilon)\|W(i, :)\|, \quad i = 1, \dots, r. \quad (79)$$

If we write $W + E_\omega = W(I + W^\dagger E_\omega)$ multiplicatively, with W^\dagger the pseudoinverse of W , then

$$W = \widehat{L}_\omega(Q'_\omega)^T(I + W^\dagger E_\omega)^{-1}.$$

Now, let $R' = (D')^{-1}R$ be the best conditioned row scaling of the triangular matrix R computed in step 1. In order to bound $\|W^\dagger E_\omega\|$, we define $Z = (D')^{-1}W$ and $E_z = (D')^{-1}E_\omega$. The equations (79) imply $\|E_z\| = O(\epsilon)\|Z\|$ and, since both D' and Z have full rank, we obtain

$$\|W^\dagger E_\omega\| = \|Z^\dagger E_z\| = O(\epsilon)\kappa(Z) = O(\epsilon\kappa(R')\kappa(Y)).$$

The last equality above is a consequence of the first equation in [7, p. 34], which implies $\|(D')^{-1}\delta W\| = O(\epsilon)\|R'\| \|Y\|$ for the error δW in the matrix multiplication of step 2 of Algorithm 5. Therefore, since $Z = R'PY^T - (D')^{-1}\delta W$, we arrive at $\kappa(Z) \leq \kappa(R')\kappa(Y)(1 + O(\epsilon)\kappa(R')\kappa(Y))$.

Thus, upon completion of step 3 of Algorithm 5, we have

$$(I + E_2)G(I + F_2) = Q' \widehat{L}_\omega(Q'_\omega)^T \quad (80)$$

with $E_2 = E'_1$, $I + F_2 = (I + F_1)(I + W^\dagger E_\omega)$ and $\|F_2\| = O(\epsilon\kappa(R')\kappa(Y))$.

Now, Theorem A.1 applied to step 4 ensures the existence of r by r matrices $\bar{U}', \bar{V}', E_L, E_R$ with \bar{U}', \bar{V}' orthogonal,

$$\begin{aligned} \|\bar{U}' - \widehat{U}_\omega\| &= O(\epsilon), & \|\bar{V}' - \widehat{V}_\omega\| &= O(\epsilon) \\ \|E_L\| &\leq O(\epsilon), & \|E_R\| &\leq O(\epsilon\kappa((D')^{-1}\widehat{L}_\omega)), \end{aligned} \quad (81)$$

and

$$\widehat{L}_\omega = (I + E_L)\bar{U}'\widehat{\Sigma}(\bar{V}')^T(I + E_R), \quad (82)$$

where $\widehat{U}_\omega\widehat{\Sigma}\widehat{V}_\omega^T$ is the SVD computed by the right Jacobi SVD algorithm on \widehat{L}_ω . Notice that we have replaced the unit row scaling of \widehat{L}_ω with the scaling given by $(D')^{-1}$. We can do this because the condition number of the former matrix is not larger than a factor \sqrt{r} times the condition number of the latter [25]. Note also that $\kappa((D')^{-1}\widehat{L}_\omega) = \kappa((D')^{-1}\widehat{L}_\omega(Q'_\omega)^T) = \kappa(Z + E_z) = \kappa(Z)(1 + O(\epsilon)\kappa(Z))$. Hence,

$$\|E_R\| = O(\epsilon\kappa(R')\kappa(Y)).$$

Substituting (82) into (80) leads to

$$(I + E_3)G(I + F_3) = Q'\bar{U}'\widehat{\Sigma}(\bar{V}')^T(Q'_\omega)^T$$

where $I + E_3 = (I + \widetilde{E}_L)^{-1}(I + E_2)$ and $I + F_3 = (I + F_2)(I + \widetilde{E}_R)^{-1}$, for $\widetilde{E}_L = Q'E_L(Q')^T$ and $\widetilde{E}_R = Q'_\omega E_R(Q'_\omega)^T$. Clearly, $\|E_3\| = O(\epsilon\kappa(X))$ and $\|F_3\| = O(\epsilon\kappa(R')\kappa(Y))$.

Finally, it only remains to show that $\widehat{U} = \mathbf{f1}(Q\widehat{U}_\omega)$ and $\widehat{V} = \mathbf{f1}(\widehat{Q}_\omega\widehat{V}_\omega)$ differ from $Q'\bar{U}'$ and $Q'_\omega\bar{V}'$ by $O(\epsilon)$. We show it for \widehat{U} , the argument for \widehat{V} is analogous. Using (78) and (81), we obtain $Q\widehat{U}_\omega = Q'\bar{U}' + O(\epsilon)$. Moreover, the standard error analysis for matrix multiplication implies that $\|\widehat{U} - Q\widehat{U}_\omega\|_F \leq nr\epsilon + O(\epsilon^2)$. The proof is concluded by observing that $\|\widehat{U} - Q'\bar{U}'\|_F \leq \|\widehat{U} - Q\widehat{U}_\omega\|_F + \|Q\widehat{U}_\omega - Q'\bar{U}'\|_F$. ■

A.1 The left-handed version

The backward error analysis above has been performed assuming that right-handed Jacobi is employed in step 4 of Algorithm 5. However, it has been observed that *left-handed* Jacobi on L_ω is usually much faster than right-handed, since the rows of L_ω are usually closer to be orthogonal than its columns (see [6, p. 565], or [21, p. 988] for a more detailed explanation of the advantages of one version over the other depending on the scaling).

The error bounds for left-handed Jacobi on an invertible matrix $A \in \mathbb{R}^{n \times n}$ remain as in Theorem A.1, at the prize of replacing the $O(\epsilon\kappa(A_N))$ by $O(\epsilon\gamma)$, where

$$\gamma = \max_{i=0,1,\dots,q} \kappa(B_i). \quad (83)$$

Here, each B_i is the diagonal scaling with unit rows of the matrix $A_i = D_i B_i$ resulting from the action of the i -th finite precision rotation along the process of left-handed Jacobi, and A_q is the first iterate satisfying the stopping criterion

$$\max_{i \neq j} \mathbf{fl} \left(\frac{|A_q(i,:)A_q(j,:)^T|}{\|A_q(i,:)\|_2 \|A_q(j,:)\|_2} \right) \leq n\epsilon, \quad \text{for } i \neq j. \quad (84)$$

To explain the origin of the additional factor γ , notice that, according to [8, Theorem 4.1], if A_i (resp. A_{i+1}) is the matrix obtained after the i -th (resp. $(i+1)$ -th) finite precision rotation, then A_{i+1} can be written as

$$A_{i+1} = R_{i+1}(A_i + \delta A_i),$$

where R_{i+1} is an exact rotation and the backward error δA_i is such that $\|\delta B_i\| \leq 72\epsilon + O(\epsilon^2)$ for the row scaling $\delta A_i = D_i \delta B_i$, where D_i is the diagonal matrix with the row norms of A_i on the diagonal. Hence,

$$A_{i+1} = R_{i+1}A_i(I + E_i)$$

with $\|E_i\| = \|A_i^{-1}\delta A_i\| = \|B_i^{-1}\delta B_i\| \leq (72\epsilon + O(\epsilon^2))\kappa(B_i)$. Notice that replacing $\|B_i^{-1}\|$ with $\kappa(B_i)$ increases the bound at most by a factor \sqrt{n} .

Repeating the argument for all q rotations up to convergence, one obtains

$$A_q = (\tilde{U}')^T A(I + \tilde{E})$$

for an exact orthogonal matrix \tilde{U}' and a matrix \tilde{E} such that $\|\tilde{E}\| \leq (72\epsilon + O(\epsilon^2))q\gamma$, with γ given by (83). The constant q in the previous error bound is pessimistic, and in fact with a finer implementation of left-handed Jacobi q can be replaced by $(s-1)p$, where s is the number of sweeps up to convergence, each of them implemented in p parallel steps [12].

Using the stopping criterion as in the end of the proof of Theorem A.1 shows that if $\hat{U}\hat{\Sigma}\hat{V}^T$ is the SVD computed by left-handed Jacobi on A with stopping criterion (84), then

$$A(I + \tilde{E}_R) = \tilde{U}'\hat{\Sigma}\tilde{V}'^T$$

for orthogonal matrices \tilde{U}' , \tilde{V}' within a distance $O(\epsilon)$ of \hat{U} , \hat{V} , and

$$\|\tilde{E}_R\| \leq 72\epsilon q\gamma + cn^2\epsilon + O(\epsilon^2) = O(\epsilon\gamma).$$

Plugging these backward errors into the proof of Theorem 2.1, we obtain for the left-handed version of Algorithm 5 (i.e., the one using left-handed Jacobi in step 4) the backward error bound

$$(I + \tilde{E})G(I + \tilde{F}) = U'\hat{\Sigma}V'^T$$

where, as in Theorem 2.1, U' and V' are orthogonal,

$$\|U' - \widehat{U}\| = O(\epsilon), \quad \|V' - \widehat{V}\| = O(\epsilon)$$

for the computed matrices \widehat{U} , $\widehat{\Sigma}$, \widehat{V} , and the backward errors satisfy

$$\|\widetilde{E}\| = O(\epsilon\kappa(X)), \quad \|\widetilde{F}\| = O(\epsilon \max\{\gamma, \kappa(R')\kappa(Y)\}),$$

with γ being the constant (83) for left-handed Jacobi on the matrix \widehat{L}_ω computed in step 3 of Algorithm 5. Therefore, the error bounds for this left-handed version of Algorithm 5 are larger than those for the right-handed one. Only if γ is of the order $O(\kappa(R')\kappa(Y))$ will high relative accuracy be achieved.

This proviso, concerning bounded growth of the condition number of matrices appearing along the Jacobi process, is analogous to the proviso, mentioned in the introduction, for the J-orthogonal algorithm. It is claimed in [8] that there is strong numerical evidence of $\gamma/\kappa(B_0) \approx 1$. A similar claim has been done in [23] on the J-orthogonal algorithm, although the sizes of the matrices in the experiments (up to 200 by 200) makes them by today's standards almost toy problems. However, the experiments in section 6 above, made for matrices of up to 1000 by 1000 for both algorithms (left-Jacobi and J-orthogonal) support this evidence as well. Hence, it seems that the increase in speed of the left-handed version is not penalized by a loss of accuracy.

B Appendix B: Proofs of Theorems 5.7 and 5.9

We begin with some previous elementary results that will be frequently used. Then, Theorems 5.7 and 5.9 will be proved in two different subsections.

Let a and a' be any two real numbers. Then

$$\frac{a - a'}{a'} = \frac{\frac{a-a'}{a}}{1 - \frac{a-a'}{a}} \quad \text{and} \quad \frac{a}{a'} = \frac{1}{1 - \frac{a-a'}{a}}. \quad (85)$$

The following Lemma bounds the relative distance between the maximum and the minimum elements in a cluster of tolerance C_l :

Lemma B.1 *Let $\Sigma_1 = \{\sigma_{i+1}, \sigma_{i+2}, \dots, \sigma_{i+d_1}\}$ be a cluster of tolerance C_l with d_1 elements. Then*

$$\frac{\sigma_{i+1} - \sigma_{i+d_1}}{\sigma_{i+1}} \leq (d_1 - 1) C_l.$$

Proof: Notice that

$$\frac{\sigma_{i+1} - \sigma_{i+d_1}}{\sigma_{i+1}} = \frac{\sigma_{i+1} - \sigma_{i+2}}{\sigma_{i+1}} + \frac{\sigma_{i+2} - \sigma_{i+3}}{\sigma_{i+1}} + \dots + \frac{\sigma_{i+d_1-1} - \sigma_{i+d_1}}{\sigma_{i+1}}.$$

Thus

$$\frac{\sigma_{i+1} - \sigma_{i+d_1}}{\sigma_{i+1}} \leq \frac{\sigma_{i+1} - \sigma_{i+2}}{\sigma_{i+1}} + \frac{\sigma_{i+2} - \sigma_{i+3}}{\sigma_{i+2}} + \dots + \frac{\sigma_{i+d_1-1} - \sigma_{i+d_1}}{\sigma_{i+d_1-1}} \leq (d_1 - 1) C_l. \quad \blacksquare$$

B.1 Proof of Theorem 5.7

Let

$$\Sigma_1 = \{\sigma_{i+1}, \sigma_{i+2}, \dots, \sigma_{i+d_1}\} \quad , \quad \Sigma_2 = \{\sigma_{i+d_1+1}, \sigma_{i+d_1+2}, \dots, \sigma_{i+d_1+d_2}\} \quad (86)$$

be the two clusters of singular values appearing in the statement of the theorem. Although in this setting the elements of Σ_1 are greater than the elements of Σ_2 , the opposite case can be proved with the notation in (86) by interchanging the roles of Σ_1 and Σ_2 .

Lemma 5.3 implies

$$rg(\Sigma_1 \cup \Sigma_2) = \min \left\{ \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}}, \frac{\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1+d_2}} \right\}, \quad (87)$$

and

$$\begin{aligned} & \min\{rg(\Sigma_1), rg(\Sigma_2)\} = \\ & = \min \left\{ \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}}, \frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1}}, \frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1+1}}, \frac{\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1+d_2}} \right\}, \end{aligned} \quad (88)$$

where if some of the subindices does not belong to $\{1, \dots, n\}$ the corresponding fraction does not appear. Therefore $rg(\Sigma_1 \cup \Sigma_2) \geq \min\{rg(\Sigma_1), rg(\Sigma_2)\}$, and the assumption (58) appearing in Theorem 5.7 leads to the following results:

1.

$$\min\{rg(\Sigma_1), rg(\Sigma_2)\} = \frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1}}. \quad (89)$$

2.

$$rg(\Sigma_1) = \frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1}}. \quad (90)$$

Thus in the setting (86), condition (58) implies that Σ_2 is the relative closest cluster to Σ_1 and it is not necessary to impose this condition explicitly. This has been done in the statement of Theorem 5.7 for the sake of clarity. Recall that one of the hypotheses of Theorem 5.7 is

$$rg(\Sigma_1) < t < 1. \quad (91)$$

The previous setting also allows to prove Theorem 5.7 in the case in which the elements of Σ_1 are smaller than the elements of Σ_2 just by interchanging the roles of Σ_1 and Σ_2 in the statement of the Theorem. Notice that condition $rg\{\Sigma_1 \cup \Sigma_2\} > \min\{rg\{\Sigma_1\}, rg\{\Sigma_2\}\}$ remains unchanged, and therefore its consequences (89), (90) still hold. This, together with $rg(\Sigma_2) < t < 1$ leads to $rg(\Sigma_1) < t$, i.e. condition (91). Therefore, in the rest of the proof we will focus on the situation (86) with assumptions (58) (and its consequences (89)-(90)) and (91).

Suppose that $(i+d_1+d_2+1) \in \{1, \dots, n\}$. If $\lambda_{\Pi(i+d_1+d_2+1)}$ is either zero or has the same sign as the elements of Λ_2 then $rg(\Lambda_2) \leq (\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1})/\sigma_{i+d_1+d_2}$. Otherwise $\lambda_{\Pi(i+d_1+d_2+1)}$ has the same sign as the elements of Λ_1 , and then $rg(\Lambda_1) \leq (\sigma_{i+d_1} - \sigma_{i+d_1+d_2+1})/\sigma_{i+d_1}$. In any case

$$\min\{rg(\Lambda_1), rg(\Lambda_2)\} \leq \max \left\{ \frac{\sigma_{i+d_1} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1}}, \frac{\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1+d_2}} \right\} = \frac{\sigma_{i+d_1} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1}}. \quad (92)$$

Suppose now that i belongs to the set $\{1, \dots, n\}$. If $\lambda_{\Pi(i)}$ has the same sign as the elements of Λ_1 then $rg(\Lambda_1) \leq (\sigma_i - \sigma_{i+1})/\sigma_{i+1}$. Otherwise $\lambda_{\Pi(i)}$ has the same sign as the elements of Λ_2 , and then $rg(\Lambda_2) \leq (\sigma_i - \sigma_{i+d_1+1})/\sigma_{i+d_1+1}$. In any case

$$\min\{rg(\Lambda_1), rg(\Lambda_2)\} \leq \max\left\{\frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}}, \frac{\sigma_i - \sigma_{i+d_1+1}}{\sigma_{i+d_1+1}}\right\} = \frac{\sigma_i - \sigma_{i+d_1+1}}{\sigma_{i+d_1+1}}. \quad (93)$$

Once (92) and (93) have been established, it only remains to prove

$$\frac{\sigma_{i+d_1} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1}} \leq R \frac{\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1+d_2}} \quad (94)$$

and

$$\frac{\sigma_i - \sigma_{i+d_1+1}}{\sigma_{i+d_1+1}} \leq R \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}}, \quad (95)$$

where

$$R = \frac{1}{1-t} \left(1 + \frac{1}{1-(d-1)C_l} + \frac{1}{1-(d-1)C_l} \frac{(d-1)C_l}{rg(\Sigma_1 \cup \Sigma_2)} \right).$$

If these two inequalities hold, then (92) and (93) imply that $\min\{rg(\Lambda_1), rg(\Lambda_2)\}$ is bounded simultaneously by the right-hand side of (94) and the right-hand side of (95). Thus, using (87) Theorem 5.7 is finally proved.

Proof of (94): Notice that:

$$\frac{\sigma_{i+d_1} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1}} = \frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1}} + \frac{\sigma_{i+d_1+1} - \sigma_{i+d_1+d_2}}{\sigma_{i+d_1}} + \frac{\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1}}. \quad (96)$$

The first term of the right-hand side in the previous equation is less than $(\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1})/\sigma_{i+d_1+d_2}$, due to (89) and (88). The third term is trivially bounded by the same quantity, since $\sigma_{i+d_1} > \sigma_{i+d_1+d_2}$. For the second term,

$$\frac{\sigma_{i+d_1+1} - \sigma_{i+d_1+d_2}}{\sigma_{i+d_1}} < \frac{\sigma_{i+d_1+1} - \sigma_{i+d_1+d_2}}{\sigma_{i+d_1+1}} \leq (d_2 - 1)C_l,$$

where the last inequality is just Lemma B.1 applied to Σ_2 . Plugging these bounds into (96) and using $rg(\Sigma_1 \cup \Sigma_2) \leq (\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1})/\sigma_{i+d_1+d_2}$, we obtain

$$\frac{\sigma_{i+d_1} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1}} \leq \left(2 + \frac{(d_2 - 1)C_l}{rg(\Sigma_1 \cup \Sigma_2)} \right) \frac{\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1+d_2}}.$$

The first factor of the right-hand side is bounded by R and (94) follows. *End of the proof of (94)* ■

Proof of (95): Notice that

$$\frac{\sigma_i - \sigma_{i+d_1+1}}{\sigma_{i+d_1+1}} = \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+d_1+1}} + \frac{\sigma_{i+1} - \sigma_{i+d_1}}{\sigma_{i+d_1+1}} + \frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1+1}}. \quad (97)$$

Now we will bound the three terms in the right-hand side of (97). We begin by the last one: using the first equality in (85), (90), (91) and (89), we get:

$$\frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1+1}} < \frac{1}{1-t} \frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1}} < \frac{1}{1-t} \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}}. \quad (98)$$

For the second term the first equality in (85) and Lemma B.1 yield

$$\frac{\sigma_{i+1} - \sigma_{i+d_1}}{\sigma_{i+d_1+1}} = \frac{\sigma_{i+d_1}}{\sigma_{i+d_1+1}} \frac{\frac{\sigma_{i+1} - \sigma_{i+d_1}}{\sigma_{i+1}}}{1 - \frac{\sigma_{i+1} - \sigma_{i+d_1}}{\sigma_{i+1}}} \leq \frac{\sigma_{i+d_1}}{\sigma_{i+d_1+1}} \frac{(d_1 - 1)C_l}{1 - (d_1 - 1)C_l}.$$

The factor $\sigma_{i+d_1}/\sigma_{i+d_1+1}$ can be bounded by $1/(1-t)$, using the second equality in (85), (90) and (91). Therefore, the following bound for the second term of the right-hand side of (97) is obtained:

$$\frac{\sigma_{i+1} - \sigma_{i+d_1}}{\sigma_{i+d_1+1}} \leq \frac{1}{1-t} \frac{(d_1 - 1)C_l}{1 - (d_1 - 1)C_l}. \quad (99)$$

Finally, the first term verifies

$$\frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+d_1+1}} = \frac{\sigma_{i+d_1}}{\sigma_{i+d_1+1}} \frac{\sigma_{i+1}}{\sigma_{i+d_1}} \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}}.$$

The factor $\sigma_{i+d_1}/\sigma_{i+d_1+1}$ has been already bounded by $1/(1-t)$, while the factor $\sigma_{i+1}/\sigma_{i+d_1}$ is bounded by $1/(1 - (d_1 - 1)C_l)$ by the second equality in (85) and Lemma B.1. Thus

$$\frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+d_1+1}} \leq \frac{1}{1-t} \frac{1}{1 - (d_1 - 1)C_l} \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}}. \quad (100)$$

Replacing (100), (99) and (98) in (97), and taking into account that $rg(\Sigma_1 \cup \Sigma_2) \leq (\sigma_i - \sigma_{i+1})/\sigma_{i+1}$,

$$\frac{\sigma_i - \sigma_{i+d_1+1}}{\sigma_{i+d_1+1}} \leq \frac{1}{1-t} \left(1 + \frac{1}{1 - (d_1 - 1)C_l} + \frac{1}{1 - (d_1 - 1)C_l} \frac{(d_1 - 1)C_l}{rg(\Sigma_1 \cup \Sigma_2)} \right) \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}}$$

is obtained. Now inequality (95) is easily proved. *End of the proof of (95).* ■

B.2 Proof of Theorem 5.9

We will only prove the case in which the elements in Σ_1 are greater than the ones in Σ_2 . In the opposite case the proof is similar, simpler and a slightly better bound can be achieved. We consider again clusters Σ_1 and Σ_2 given by (86). Lemma 5.3 and the fact that Σ_2 is the relative closest cluster to Σ_1 imply

$$rg(\Sigma_1) = \frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1}} < t. \quad (101)$$

We will split the proof in two cases, depending on which of the values in (87) equals $rg(\Sigma_1 \cup \Sigma_2)$.

Case 1: $rg(\Sigma_1 \cup \Sigma_2) = \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}}$

Assumption (59), together with (101) and Lemma 5.3, yield in this case

$$\frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}} = rg(\Sigma_1 \cup \Sigma_2) = \min\{rg(\Sigma_1), rg(\Sigma_2)\} \leq rg(\Sigma_1) = \frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1}} \leq \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}}.$$

Therefore

$$rg(\Sigma_1) = \frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1}} = \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}}. \quad (102)$$

If $\lambda_{\Pi(i)}$ has the same sign as the elements in Λ_1 then:

$$rg(\Lambda_1) \leq \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}} = rg(\Sigma_1).$$

Using a trivial lower (resp. upper) bound for $rg(\Lambda_1)$ (resp. $rg(\Sigma_1)$) one immediately obtains the bound in Theorem 5.9. Otherwise $\lambda_{\Pi(i)}$ has the same sign as the elements in Λ_2 , so

$$rg(\Lambda_2) \leq \frac{\sigma_i - \sigma_{i+d_1+1}}{\sigma_{i+d_1+1}} = \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+d_1+1}} + \frac{\sigma_{i+1} - \sigma_{i+d_1}}{\sigma_{i+d_1+1}} + \frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1+1}}. \quad (103)$$

Now the three terms in the right-hand side of (103) will be bounded. We begin by the last one. Using the first equality in (85) and (101):

$$\frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1+1}} < \frac{1}{1-t} \frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1}} = \frac{1}{1-t} rg(\Sigma_1) \quad (104)$$

is obtained. For the second term, we get again (99) following exactly the same steps (notice that (90) and (91) are just (101)). The first term in the right-hand side of (103) verifies:

$$\frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+d_1+1}} = \frac{\sigma_{i+d_1}}{\sigma_{i+d_1+1}} \frac{\sigma_{i+1}}{\sigma_{i+d_1}} \frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+1}} = \frac{\sigma_{i+d_1}}{\sigma_{i+d_1+1}} \frac{\sigma_{i+1}}{\sigma_{i+d_1}} rg(\Sigma_1),$$

where (102) has been used. Thus the same argument used to get (100) leads to

$$\frac{\sigma_i - \sigma_{i+1}}{\sigma_{i+d_1+1}} \leq \frac{1}{1-t} \frac{1}{1 - (d_1 - 1)C_l} rg(\Sigma_1). \quad (105)$$

Replacing (105), (99) and (104) in (103),

$$rg(\Lambda_2) \leq \frac{1}{1-t} \left(1 + \frac{1}{1 - (d_1 - 1)C_l} + \frac{1}{1 - (d_1 - 1)C_l} \frac{(d_1 - 1)C_l}{rg(\Sigma_1)} \right) rg(\Sigma_1)$$

is obtained. Using a lower (resp. upper) bound for the left-hand side (resp. right-hand side) of the previous inequality, one immediately gets the bound in Theorem 5.9.

Case 2: $rg(\Sigma_1 \cup \Sigma_2) = \frac{\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1+d_2}}$

Assumption (59) and Lemma 5.3, yield in this case

$$\frac{\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1+d_2}} = rg(\Sigma_1 \cup \Sigma_2) = \min\{rg(\Sigma_1), rg(\Sigma_2)\} \leq rg(\Sigma_2) \leq \frac{\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1+d_2}}.$$

Therefore,

$$rg(\Sigma_2) = \frac{\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1+d_2}} \leq rg(\Sigma_1). \quad (106)$$

If $\lambda_{\Pi(i+d_1+d_2+1)}$ is zero or has the same sign as the elements in Λ_2 then

$$rg(\Lambda_2) \leq \frac{\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1+d_2}} = rg(\Sigma_2) \leq rg(\Sigma_1).$$

Using a trivial lower (resp. upper) bound for $rg(\Lambda_2)$ (resp. $rg(\Sigma_1)$) one immediately obtains the bound in Theorem 5.9. Otherwise $\lambda_{\Pi(i+d_1+d_2+1)}$ has the same sign as the elements in Λ_1 , so

$$rg(\Lambda_1) \leq \frac{\sigma_{i+d_1} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1}} = \frac{\sigma_{i+d_1} - \sigma_{i+d_1+1}}{\sigma_{i+d_1}} + \frac{\sigma_{i+d_1+1} - \sigma_{i+d_1+d_2}}{\sigma_{i+d_1}} + \frac{\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1}}. \quad (107)$$

Now we bound the three terms in the right-hand side of (107). The first one is just $rg(\Sigma_1)$ by (101). Lemma B.1 applied to Σ_2 leads to the bound

$$\frac{\sigma_{i+d_1+1} - \sigma_{i+d_1+d_2}}{\sigma_{i+d_1}} < \frac{\sigma_{i+d_1+1} - \sigma_{i+d_1+d_2}}{\sigma_{i+d_1+1}} \leq (d_2 - 1)C_l$$

for the second term. The third term verifies

$$\frac{\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1}} < \frac{\sigma_{i+d_1+d_2} - \sigma_{i+d_1+d_2+1}}{\sigma_{i+d_1+d_2}} = rg(\Sigma_2) \leq rg(\Sigma_1)$$

due to (106). Then one obtains from (107)

$$rg(\Lambda_1) \leq \left(2 + \frac{(d_2 - 1)C_l}{rg(\Sigma_1)} \right) rg(\Sigma_1).$$

Using a lower (resp. upper) bound for the left-hand side (resp. right-hand side) of the previous inequality, one immediately gets the bound in Theorem 5.9.

Acknowledgements. The authors thank Prof. Zlatko Drmač, who provided the source code for the one-sided Jacobi SVD routine employed in the experiments. As can be seen in the numerical tests in Section 6, the performance of his code is excellent. The authors thank also Prof. J. W. Demmel for providing the source code of the routines used in [6].

References

- [1] E. ANDERSON, Z. BAI, C. BISCHOF, S. BLACKFORD, J. DEMMEL, J. DONGARRA, J. DU CROZ, A. GREENBAUM, S. HAMMARLING, A. MCKENNEY AND D. SORESENSEN, *LAPACK User's Guide*, 3rd Ed., SIAM, Philadelphia, 1999.
- [2] J. BARLOW AND J. DEMMEL, *Computing accurate eigensystems of scaled diagonally dominant matrices*, SIAM J. Numer. Anal., 27 (1990), pp. 762-791.
- [3] J.R. BUNCH AND B. PARLETT, *Direct methods for solving symmetric indefinite systems of linear equations*, SIAM J. Numer. Anal., 8 (1971), pp. 639-655.
- [4] A.J. COX AND N.J. HIGHAM, *Stability of Householder QR factorization for weighted least squares problems*, in Numerical Analysis 1997 (Dundee), Proceedings of the 17th Dundee Biennial Conference, Pitman Research Notes in Mathematics vol. 380, pp. 57-73, Longman, UK, 1998.
- [5] C. DAVIS AND W. KAHAN, *The rotation of eigenvectors by a perturbation III*, SIAM J. Numer. Anal., 7 (1970), pp. 1-46.
- [6] J. DEMMEL, *Accurate singular value decompositions of structured matrices*, SIAM J. Matrix Anal. Appl., 21 (1999), pp. 562-580.
- [7] J. DEMMEL, M. GU, S. EISENSTAT, I. SLAPNIČAR, K. VESELIĆ AND Z. DRMAČ, *Computing the singular value decomposition with high relative accuracy*, Linear Algebra Appl. 299 (1999), pp. 21-80.
- [8] J. DEMMEL AND K. VESELIĆ, *Jacobi's method is more accurate than QR*, SIAM J. Matrix Anal. Appl., 13 (1992), pp. 1204-1245.

- [9] F. M. DOPICO, *A note on $\sin \Theta$ theorems for singular subspace variations*, BIT - Numerical Mathematics, 40 (2000), pp. 395-403.
- [10] F. M. DOPICO AND J. MORO, *Perturbation theory for simultaneous bases of singular subspaces*, BIT - Numerical Mathematics, 42 (2002), pp. 84-109.
- [11] Z. DRMAČ, *Implementation of Jacobi rotations for accurate singular value computation in floating-point arithmetic*, SIAM J. Sci. Comput., 18 (1997), pp. 1200-1222.
- [12] Z. DRMAČ, *Accurate computation of the product-induced singular value decomposition with applications*, SIAM J. Numer. Anal., 35 (1998), pp. 1969-1994.
- [13] Z. DRMAČ, *A posteriori computation of the singular vectors in a preconditioned Jacobi SVD algorithm*, IMA J. Numer. Anal., 19 (1999), pp. 191-213.
- [14] Z. DRMAČ AND K. VESELIĆ, *Approximate eigenvectors as preconditioner*, Linear Algebra Appl. 309 (2000), pp. 191-215.
- [15] S. EISENSTAT AND I. IPSEN, *Relative perturbation techniques for singular value problems*, SIAM J. Numer. Anal., 32 (1995), pp. 1972-1988.
- [16] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, 3rd. ed., The Johns Hopkins University Press, Baltimore, 1996.
- [17] N. J. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, SIAM, Philadelphia, 1996.
- [18] N. J. HIGHAM, *QR factorization with complete pivoting and accurate computation of the SVD*, Linear Algebra Appl., 309 (2000), pp. 153-174.
- [19] R.-C. LI, *Relative perturbation theory: (I) Eigenvalue and singular value variations*, SIAM J. Matrix Anal. Appl., 19 (1998), pp. 956-982.
- [20] R.-C. LI, *Relative Perturbation Theory: (II) Eigenspace and Singular Subspace Variations*, SIAM J. Matrix Anal. Appl., 20 (1999), pp. 471-492.
- [21] R. MATHIAS, *Accurate eigensystem computations by Jacobi methods*, SIAM J. Matrix Anal. Appl., 16 (1995), pp. 977-1003.
- [22] M. POWELL AND J. REID, *On applying Householder transformations to linear least squares problems*, in Proceedings of the IFIP Congress 1968, pp. 122-126, North-Holland, Amsterdam, 1969.
- [23] I. SLAPNIČAR, *Accurate Symmetric Eigenreduction by a Jacobi Method*, PhD Thesis, Fachbereich Mathematik Fernuniversität, Gesamthochschule Hagen, Germany, 1992.
- [24] I. SLAPNIČAR, *Componentwise analysis of direct factorization of real symmetric and Hermitian matrices*, Linear Algebra Appl., 272 (1998), pp. 227-275.
- [25] A. VAN DER SLUIS, *Condition numbers and equilibration of matrices*, Numer. Math., 14 (1969), pp. 14-23.
- [26] G. W. STEWART AND J.-G. SUN, *Matrix Perturbation Theory*, Academic Press, New York, 1990.

- [27] K. VESELIĆ, *A Jacobi eigenreduction algorithm for definite matrix pairs*, Numer. Math., 64 (1993), pp. 241-269.
- [28] K. VESELIĆ AND I. SLAPNIČAR, *Floating-point perturbations of Hermitian matrices*, Linear Algebra Appl., 195 (1993), pp. 81-116.