

On the Solution of Wahba's Problem on $SO(n)$

Anton H. J. de Ruiter · James Richard
Forbes

Abstract Wahba's problem is an important problem in both aerospace and robotics fields. It typically involves finding an optimal rotation to fit a series of vector measurements. In this paper, a rigorous analysis of the famous Wahba problem on $SO(n)$ is presented. The entire set of solutions is obtained using both a singular value decomposition and matrix square-root method. Conditions for uniqueness of solutions are also obtained, correcting some errors in the existing literature. It is then shown that under a mild condition, Wahba's problem can be recast as a convex optimization problem with linear matrix inequality constraints. This opens the door to a whole host of new possible solution methods for these types of Wahba problems.

Keywords Wahba's Problem · $SO(n)$ · Matrix Square Root · Singular Value Decomposition · Linear Matrix Inequalities

Introduction

Wahba's problem is an important problem in engineering, which typically involves finding an optimal rotation to fit a series of vector measurements. It has received significant attention in both the aerospace and robotics fields,

A. de Ruiter
Department of Aerospace Engineering, Ryerson University
Toronto, Ontario, Canada
Tel.: 416-979-5000 ext. 4878
Fax: 416-979-5056
E-mail: aderuiter@ryerson.ca

J. Forbes
Department of Aerospace Engineering, University of Michigan
1320 Beal Avenue, Ann Arbor, Michigan, USA, 48109.
Tel.: 734-763-5214
Fax: 734-763-0578
E-mail: forbesrj@umich.edu

seemingly independently of each other. In particular, the name “Wahba’s Problem” stems from the aerospace field, and is not typically used in the robotics field. In the aerospace field, Wahba’s problem was introduced in 1965 by Grace Wahba [1], and attention has been focused on solving Wahba’s problem on the special orthogonal group $SO(3)$, which consists of all rotation matrices transforming three-dimensional right-handed reference frames. In the robotics field, the more general problem on $SO(n)$ has been considered. Restricting attention to $SO(3)$ has allowed researchers to develop solutions using a unit quaternion parameterization of $SO(3)$. In the realm of aerospace engineering the original solution using quaternions is attributed to Davenport [2], while in the robotics realm the original solution using quaternions is attributed to Horn [3]. It should be noted that Horn considers both translation and rotation, while Davenport considers only rotation. There have been many different methods developed to solve Wahba’s problem on $SO(3)$, primarily in the aerospace field, the most famous one being QUEST [4]. A good survey of the different methods may be found in [5], and the references therein. However, since most of these methods are based on the unit quaternion, they do not have any applicability beyond $SO(3)$. Solutions based on the singular value decomposition and the matrix square-root, however, are directly applicable to $SO(n)$, and have received attention by several authors ([6] to [12]). Recently, a useful generalization of Wahba’s problem on $SO(3)$ has been made, allowing the determination of both attitude and body-rate using a time-history of vector measurements (see [13] to [15]). However, this paper does not examine this problem, and treats only Wahba’s original problem.

In this paper, a thorough, self-contained and rigorous analysis of Wahba’s problem on $SO(n)$ is provided, based on the singular value decomposition. The matrix square-root solution is also presented as a simple corollary of the singular value decomposition solution. While the forms of the singular value decomposition and the matrix square-root solutions have been previously presented in the literature, there are aspects that have not received thorough attention. In particular, a distinguishing feature of the treatment in this paper is that the entire set of solutions is derived. This has not received much attention in the existing literature for Wahba’s problem on $SO(n)$. Using this set, necessary and sufficient conditions for uniqueness of solutions are obtained. In doing so, previously presented necessary and sufficient conditions for uniqueness of solutions are corrected. Finally, having presented a thorough analysis of the solutions to Wahba’s problem, the results from [16] are generalized to $SO(n)$, by showing that under a mild condition (that is satisfied in many practical applications), Wahba’s problem on $SO(n)$ may be recast as a linear matrix inequality (LMI) optimization problem. This opens the door to a whole new class of solvers for Wahba’s problem on $SO(n)$. In addition, it suggests an approach for how non-convex $SO(n)$ constraints may potentially be relaxed to convex LMI constraints, in other optimization problems involving elements of $SO(n)$, making them more tractable.

The remainder of the paper is organized as follows. First, some notations used in this paper are explained. After this, Wahba’s problem and two related

problems are defined, demonstrating that they have the same solution. Following this, a rigorous derivation and analysis of solutions is presented using the singular value decomposition methods. From these results, the matrix square-root solution is then obtained. Next, it is demonstrated that under a mild condition, the Wahba problem may be recast as an LMI problem, leading to an identical solution. The paper concludes with a numerical example, followed by concluding remarks. The appendix contains mathematical facts, which are used throughout the paper.

Notation

We use the following notation for the gradient of a scalar function of a matrix (of which a scalar function of a column vector is a special case). Let $\mathbf{X} \in \mathbb{R}^{n \times m}$, and $f(\mathbf{X}) : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$. We then define the derivative of f with respect to the matrix \mathbf{X} , as

$$\frac{\partial f}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial f}{\partial X_{11}} & \cdots & \frac{\partial f}{\partial X_{1m}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial X_{n1}} & \cdots & \frac{\partial f}{\partial X_{nm}} \end{bmatrix}, \quad (1)$$

where X_{ij} is the ij^{th} term of \mathbf{X} . Note that the chosen derivative convention in (1) is one of two possibilities for the Jacobian of $f(\mathbf{X})$, the other being its transpose.

We use the following notation for a block-diagonal matrix

$$\text{diag}_{i=1, \dots, n} \{\mathbf{X}_i\} = \text{diag}\{\mathbf{X}_1, \dots, \mathbf{X}_n\} = \begin{bmatrix} \mathbf{X}_1 & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \mathbf{X}_n \end{bmatrix}, \quad (2)$$

where \mathbf{X}_i may be scalars or square matrices. The trace of a square matrix \mathbf{X} will be denoted by $\text{tr}[\mathbf{X}]$, and $\mathbf{1}_{n \times n}$ denotes an n by n identity matrix. Finally, $\|\mathbf{A}\|$ will denote the induced 2-norm of the matrix \mathbf{A} .

Problem Statements

In this section, we shall define three different problems, all of which have the same solution, namely the solution to Wahba's problem. We start by defining the original Wahba problem on $SO(n)$.

Problem 1 (Wahba's Problem) *Given N vectors, $\mathbf{s}_{b,k} \in \mathbb{R}^n$, with corresponding vectors $\mathbf{s}_{I,k} \in \mathbb{R}^n$, Wahba's Problem is to find the matrix $\mathbf{C} \in SO(n)$, where $SO(n) = \{\mathbf{C} \in \mathbb{R}^{n \times n} : \mathbf{C}^T \mathbf{C} = \mathbf{C} \mathbf{C}^T = \mathbf{1}, \det \mathbf{C} = +1\}$, to minimize the cost function*

$$J = \sum_{k=1}^N w_k (\mathbf{s}_{b,k} - \mathbf{C} \mathbf{s}_{I,k})^T (\mathbf{s}_{b,k} - \mathbf{C} \mathbf{s}_{I,k}), \quad (3)$$

where $0 < w_k < \infty$ are positive weights for $k = 1, \dots, N$.

This is a generalization to $SO(n)$ of the well-known problem originally posed by Grace Wahba in 1965 [1], where $\mathbf{C} \in SO(3)$ is a rotation matrix describing the transformation between right-handed reference frames.

Expanding (3), we obtain

$$J = \sum_{k=1}^N w_k \mathbf{s}_{b,k}^T \mathbf{s}_{b,k} + \sum_{k=1}^N w_k \mathbf{s}_{I,k}^T \mathbf{C}^T \mathbf{C} \mathbf{s}_{I,k} - 2 \sum_{k=1}^N w_k \mathbf{s}_{b,k}^T \mathbf{C} \mathbf{s}_{I,k}.$$

Noting that $\mathbf{C}^T \mathbf{C} = \mathbf{1}$ for all $\mathbf{C} \in SO(n)$, the cost function becomes

$$J = \sum_{k=1}^N w_k \mathbf{s}_{b,k}^T \mathbf{s}_{b,k} + \sum_{k=1}^N w_k \mathbf{s}_{I,k}^T \mathbf{s}_{I,k} - 2 \text{tr} [\mathbf{C} \mathbf{B}^T], \quad \forall \mathbf{C} \in SO(n), \quad (4)$$

where

$$\mathbf{B}^T = \sum_{k=1}^N w_k \mathbf{s}_{I,k} \mathbf{s}_{b,k}^T. \quad (5)$$

Clearly, only the third term in (4) depends on \mathbf{C} . Therefore minimizing J in (3) subject to $\mathbf{C} \in SO(n)$ is equivalent to solving the maximization problem

$$\text{maximize } \hat{J} = \text{tr} [\mathbf{C} \mathbf{B}^T] \text{ subject to } \mathbf{C} \in SO(n). \quad (6)$$

In some instances, it may be desirable to orthonormalize a given matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$. For example, when $\mathbf{C} \in SO(3)$, its kinematics satisfy Poisson's equation $\dot{\mathbf{C}} = -\boldsymbol{\omega}^\times \mathbf{C}$ [17], with initial condition $\mathbf{C}(t_0) \in SO(3)$. Direct numerical integration will result in a solution $\hat{\mathbf{C}}(t)$ which is no longer in $SO(3)$. It is therefore desirable to orthonormalize $\hat{\mathbf{C}}(t)$ after each numerical integration step. Orthonormalization of a rotation matrix estimate $\hat{\mathbf{C}}(t)$ can be thought of in the same manner as normalization of a quaternion estimate.

Considering the matrix \mathbf{D} to be an approximation of a matrix $\mathbf{C} \in SO(n)$, it is reasonable to expect $\det[\mathbf{D}] > 0$ (otherwise it would be a very poor approximation). It would be desirable to orthonormalize the matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$ in an optimal manner, as stated in the next Problem. Note that this is very similar to the orthogonal Procrustes problem [7], in which \mathbf{C} is only required to be orthonormal, without any restriction on the sign of its determinant.

Problem 2 (Matrix Orthonormalization) Let $\mathbf{D} \in \mathbb{R}^{n \times n}$,

$$\text{minimize } J = \text{tr} [(\mathbf{D} - \mathbf{C})^T (\mathbf{D} - \mathbf{C})], \quad (7)$$

subject to $\mathbf{C} \in SO(n)$.

As will be shown later (see the discussion following Theorem 3), the solution of Problem 2 is equivalent to the solution of the orthogonal Procrustes problem when $\det \mathbf{D} > 0$.

Analogously to Wahba's problem (Problem 1), it is readily shown that the minimization problem in (7) is equivalent to the maximization problem

$$\text{maximize } \hat{J} = \text{tr} [\mathbf{C}\mathbf{D}^T] \text{ subject to } \mathbf{C} \in SO(n). \quad (8)$$

The previous two problems are purely rotational, and have received significant attention in the aerospace community, particularly for the purpose of spacecraft attitude determination. The final problem considered in this paper is that of simultaneous translation and rotation determination, which is a problem on $SE(n) = SO(n) \times \mathbb{R}^n$. In particular, the problem is given two point clouds obtained from a vehicle's vision (or other) system at two separate time instants, determine the translation and rotation of the vehicle during the time interval. This problem has received significant attention in the robotics community.

Problem 3 (Simultaneous Translation and Rotation) *Given N vectors, $\mathbf{p}_{b,k} \in \mathbb{R}^n$, with corresponding vectors $\mathbf{p}_{I,k} \in \mathbb{R}^n$, find $(\mathbf{C}, \mathbf{r}) \in SE(n)$ to minimize the cost function*

$$J = \sum_{k=1}^N w_k \left(\mathbf{p}_{b,k} - \mathbf{C}(\mathbf{p}_{I,k} - \mathbf{r}) \right)^T \left(\mathbf{p}_{b,k} - \mathbf{C}(\mathbf{p}_{I,k} - \mathbf{r}) \right), \quad (9)$$

where $0 < w_k < \infty$ are positive weights for $k = 1, \dots, N$.

It is well-known that this can be reduced to a minimization problem in \mathbf{C} [3, 11]. Indeed, consider the change of variables

$$\mathbf{r} = -\mathbf{C}^T \mathbf{p}_b + \mathbf{p}_I + \delta \mathbf{r}, \quad (10)$$

where

$$\mathbf{p}_b = \frac{1}{w} \sum_{k=1}^N w_k \mathbf{p}_{b,k}, \quad \mathbf{p}_I = \frac{1}{w} \sum_{k=1}^N w_k \mathbf{p}_{I,k}, \quad w = \sum_{k=1}^N w_k. \quad (11)$$

Substituting this into the cost function (9), and enforcing the constraint $\mathbf{C} \in SO(n)$, the cost function becomes

$$\begin{aligned} J &= \sum_{k=1}^N w_k \left((\mathbf{p}_{b,k} - \mathbf{p}_b) - \mathbf{C}(\mathbf{p}_{I,k} - \mathbf{p}_I) \right)^T \left((\mathbf{p}_{b,k} - \mathbf{p}_b) - \mathbf{C}(\mathbf{p}_{I,k} - \mathbf{p}_I) \right) \\ &\quad + 2 \sum_{k=1}^N w_k \left((\mathbf{p}_{b,k} - \mathbf{p}_b) - \mathbf{C}(\mathbf{p}_{I,k} - \mathbf{p}_I) \right)^T \mathbf{C} \delta \mathbf{r} \\ &\quad + \sum_{k=1}^N w_k \delta \mathbf{r}^T \delta \mathbf{r}. \end{aligned} \quad (12)$$

By the definition of \mathbf{p}_b and \mathbf{p}_I in (11), we find that the second term in (12) vanishes, and Problem 3 is solved by $\delta \mathbf{r} = \mathbf{0}$, together with $\mathbf{C} \in SO(n)$ minimizing

$$J_1 = \sum_{k=1}^N w_k \left((\mathbf{p}_{b,k} - \mathbf{p}_b) - \mathbf{C}(\mathbf{p}_{I,k} - \mathbf{p}_I) \right)^T \left((\mathbf{p}_{b,k} - \mathbf{p}_b) - \mathbf{C}(\mathbf{p}_{I,k} - \mathbf{p}_I) \right), \quad (13)$$

The minimization of (13) is identical to the solution of Problem 1. With $\delta \mathbf{r} = \mathbf{0}$, equation (10) yields the solution for \mathbf{r}

$$\mathbf{r} = -\mathbf{C}^T \mathbf{p}_b + \mathbf{p}_I, \quad (14)$$

where \mathbf{C} is the minimizing solution of (13). Finally, as for Problems 1 and 2, it can be readily shown that the minimization of J_1 in (13), is equivalent to the maximization problem

$$\text{maximize } \hat{J}_1 = \text{tr} [\mathbf{C} \bar{\mathbf{B}}^T] \text{ subject to } \mathbf{C} \in SO(n), \quad (15)$$

where

$$\bar{\mathbf{B}}^T = \sum_{k=1}^N w_k (\mathbf{p}_{I,k} - \mathbf{p}_I)(\mathbf{p}_{b,k} - \mathbf{p}_b)^T. \quad (16)$$

Comparing (6), (8), and (15), it can be seen that Problems 1, 2 and 3 have identical solutions. As such, we shall only consider Problem 1 from this point on.

Problem 1 Solution

We first examine a solution to Problem 1 based on the singular value decomposition. This has been previously considered in references [9,11] for problems in $SO(n)$, and in [10,12] for problems in $SO(3)$. We shall fully develop the solution, and examine conditions for uniqueness of solutions, correcting some erroneous uniqueness conditions in the literature.

Consider a singular value decomposition of \mathbf{B} in (6), given by

$$\mathbf{B} = \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T, \quad (17)$$

where $\mathbf{V}^T \mathbf{V} = \mathbf{1}$, $\mathbf{U}^T \mathbf{U} = \mathbf{1}$, $\boldsymbol{\Sigma} = \text{diag}_{i=1, \dots, m} \{\sigma_i \mathbf{1}_{n_i \times n_i}\}$ and $\sigma_1 > \dots > \sigma_m \geq 0$ are the distinct singular values of \mathbf{B} , with multiplicity n_i , such that $\sum_{i=1}^m n_i = n$. Next, fixing \mathbf{V} and \mathbf{U} in (17) (as shown in Lemma 6 in the appendix, they are not unique), consider the change of variables given by

$$\mathbf{S} = \mathbf{V}^T \mathbf{C} \mathbf{U}. \quad (18)$$

Defining the set

$$E = \{\mathbf{S} \in \mathbb{R}^{n \times n} : \mathbf{S}^T \mathbf{S} = \mathbf{S} \mathbf{S}^T = \mathbf{1}, \det \mathbf{S} = \det \mathbf{U} \det \mathbf{V}\}, \quad (19)$$

it is readily verified that the change of variables in (18) defines a bijective mapping from $SO(n)$ to E , with the inverse mapping given by

$$\mathbf{C} = \mathbf{V} \mathbf{S} \mathbf{U}^T. \quad (20)$$

Therefore, substituting (17) and (20) into (6), the maximization problem in (6) (which is equivalent to Problem 1) is equivalent to

$$\text{maximize } \hat{J} = \text{tr} [\mathbf{S} \boldsymbol{\Sigma}] \text{ subject to } \mathbf{S} \in E, \quad (21)$$

where E is given by (19). Since the mapping in (18) is bijective, solutions of (6) (and consequently of Problem 1) are unique if and only if solutions of (21) are unique. Furthermore, identification of the entire set of solutions to (21), leads through (20) to the entire set of solutions to Problem 1.

Noting that $|\det \mathbf{U} \det \mathbf{V}| = 1$, let us now relax the constraint on \mathbf{S} , to simply $\mathbf{S}^T \mathbf{S} = \mathbf{1}$. That is, we allow both $\det \mathbf{S} = \pm 1$. Consequently, we now consider the maximization problem

$$\text{maximize } \hat{J} = \text{tr} [\mathbf{S}\boldsymbol{\Sigma}] \text{ subject to } \mathbf{S} \in O(n), \quad (22)$$

where $O(n) = \{\mathbf{S} \in \mathbb{R}^{n \times n} : \mathbf{S}^T \mathbf{S} = \mathbf{S}\mathbf{S}^T = \mathbf{1}\}$. Since the constraint set $\mathbf{S} \in O(n)$ is compact, a global maximizing solution of (22) exists. Clearly, if a maximizing solution of (22) satisfies $\det \mathbf{S} = \det \mathbf{U} \det \mathbf{V}$, then it is also a maximizing solution of (21). Since E and

$$O(n) - E = \{\mathbf{S} \in \mathbb{R}^{n \times n} : \mathbf{S}^T \mathbf{S} = \mathbf{S}\mathbf{S}^T = \mathbf{1}, \det \mathbf{S} = -\det \mathbf{U} \det \mathbf{V}\},$$

are disjoint and compact sets, a global maximizing solution of (21) also exists and will be a local maximizing solution of (22). Therefore, we can use the necessary conditions for (22) to identify global maximizing solutions of (21).

We note that the constraint $\mathbf{S}^T \mathbf{S} = \mathbf{1}$ is equivalent to $\mathbf{S}\mathbf{S}^T = \mathbf{1}$. Note that the matrix constraint $\mathbf{S}\mathbf{S}^T = \mathbf{1}$ is symmetric, and it therefore contains $n(n+1)/2$ independent scalar constraint equations, contained in the upper-triangular part, which are given by

$$\phi_{ij}(\mathbf{S}) = 0, \quad i = 1, \dots, n, \quad j = i, \dots, n, \quad (23)$$

where

$$\phi_{ij}(\mathbf{S}) = \mathbf{e}_i^T (\mathbf{1} - \mathbf{S}\mathbf{S}^T) \mathbf{e}_j \quad i = 1, \dots, n, \quad j = i, \dots, n,$$

and $\mathbf{e}_i \in \mathbb{R}^n$ has the i^{th} equal to one, and all other entries equal to zero.

It is clear that the performance index \hat{J} and the constraint in (22) are smooth. We shall now verify that the required constraint qualification (linear independence of the constraint gradients) is also satisfied, which allows the use of the Kuhn-Tucker necessary conditions to identify all local maximizing solutions [18, p. 328]. From (23), it is readily found that the constraint gradients are given by

$$\frac{\partial \phi_{ij}}{\partial \mathbf{S}} = -\mathbf{X}_{ij} \mathbf{S}, \quad i = 1, \dots, n, \quad j = i, \dots, n. \quad (24)$$

where

$$\mathbf{X}_{ij} = [\mathbf{e}_i \mathbf{e}_j^T + \mathbf{e}_j \mathbf{e}_i^T], \quad i = 1, \dots, n, \quad j = i, \dots, n. \quad (25)$$

The matrices \mathbf{X}_{ii} for $i = 1, \dots, n$ in (25) contain a two in the i^{th} entry, with all other entries equal to zero. The matrices \mathbf{X}_{ij} for $i = 1, \dots, n, j = i+1, \dots, n$ contain ones in the ij^{th} and ji^{th} entries, with all other entries are equal to zero. The collection of matrices \mathbf{X}_{ij} in (25) are therefore linearly independent. Since \mathbf{S} is nonsingular (it is in $O(n)$), it follows that the constraint gradients in (24) are linearly independent for all $\mathbf{S} \in O(n)$ (including the local maximizing

solutions). We now conclude that any local maximizing solution \mathbf{S} of (22) must satisfy [18, p. 328]

$$\frac{\partial \mathcal{L}}{\partial \mathbf{S}} = \mathbf{0}, \quad \mathbf{S} \in O(n), \quad (26)$$

where

$$\mathcal{L} = \text{tr} [\mathbf{S}\boldsymbol{\Sigma}] + \text{tr} [\boldsymbol{\Gamma} (\mathbf{1} - \mathbf{S}\mathbf{S}^T)], \quad (27)$$

is the Lagrangian, and $\boldsymbol{\Gamma} \in \mathbb{R}^{n \times n}$ is a matrix of Lagrange multipliers. Note that $\boldsymbol{\Gamma}$ lower triangular, since all independent constraints are contained in the upper-triangular part of $\mathbf{1} - \mathbf{S}\mathbf{S}^T = \mathbf{0}$.

Note that linear independence of the constraint gradients is precisely why we dropped the determinant requirement in (22), since $\det \mathbf{S} = \det \mathbf{U} \det \mathbf{V}$ is not independent of the orthonormality condition $\mathbf{S}^T \mathbf{S} = \mathbf{1}$. In fact, it can readily be shown that gradient of the determinant constraint is linearly dependent on the gradients of the orthonormality constraints given in (24). This technicality has been overlooked in [11].

The necessary condition in (26) now becomes

$$\boldsymbol{\Sigma} - (\boldsymbol{\Gamma} + \boldsymbol{\Gamma}^T) \mathbf{S} = \mathbf{0}, \quad \mathbf{S}\mathbf{S}^T = \mathbf{1} \quad (28)$$

Denoting $\mathbf{W} = \boldsymbol{\Gamma} + \boldsymbol{\Gamma}^T$, we find that the necessary condition for a solution to (22) is that \mathbf{S} must satisfy

$$\mathbf{W}\mathbf{S} = \boldsymbol{\Sigma}, \quad (29)$$

where \mathbf{W} is a real symmetric matrix. Multiplying (29) by its transpose, and enforcing the constraint $\mathbf{S}\mathbf{S}^T = \mathbf{1}$, we find that

$$\mathbf{W}^2 = \boldsymbol{\Sigma}^2. \quad (30)$$

As such, the matrix \mathbf{W} is a real symmetric matrix square root of the diagonal positive semi-definite matrix $\boldsymbol{\Sigma}^2$. Now, there are many possible choices for \mathbf{W} , and we must find the one(s) which maximize the performance index \hat{J} in (21), while yielding $\det \mathbf{S} = \det \mathbf{U} \det \mathbf{V}$. To do this, we first identify all symmetric square-roots of $\boldsymbol{\Sigma}^2$. Having done this, we then identify associated $\mathbf{S} \in E$ (see (19)) which satisfy the necessary condition (29), and select the one(s) that maximize \hat{J} in (21).

Lemma 1 *The set of all real symmetric \mathbf{W} satisfying (30) is given by*

$$\mathbf{W} = \boldsymbol{\Delta} \text{diag} \{ \bar{\boldsymbol{\Sigma}}_i \} \boldsymbol{\Delta}^T, \quad (31)$$

where

$$\bar{\boldsymbol{\Sigma}}_i = \text{diag} \{ p_i^j \sigma_i \}, \quad p_i^j = \pm 1, \quad j = 1, \dots, n_i, \quad (32)$$

for $i = 1, \dots, m$, and

$$\boldsymbol{\Delta} = \text{diag} \{ \boldsymbol{\Delta}_i \}, \quad (33)$$

where $\boldsymbol{\Delta}_i \in \mathbb{R}^{n_i \times n_i}$ are real orthonormal matrices.

Proof First, consider any \mathbf{W} of the form in (31). Then, direct calculation yields $\mathbf{W}^2 = \boldsymbol{\Sigma}^2$.

Conversely, suppose that \mathbf{W} is real and symmetric, and satisfies (30). Then, by Lemma 4 (see the appendix) it can be written in the form

$$\mathbf{W} = \bar{\mathbf{V}}\bar{\mathbf{A}}\bar{\mathbf{V}}^T, \quad (34)$$

where $\bar{\mathbf{V}} \in \mathbb{R}^{n \times n}$ is a real orthonormal matrix, and

$$\bar{\mathbf{A}} = \text{diag} \{ \bar{\lambda}_i \}, \quad (35)$$

$i=1, \dots, n$

is a real diagonal matrix containing the eigenvalues of \mathbf{W} . The columns of $\bar{\mathbf{V}}$ contain eigenvectors of \mathbf{W} corresponding to each of the eigenvalues in $\bar{\mathbf{A}}$. Substituting (34) into (30), we obtain

$$\boldsymbol{\Sigma}^2 = \bar{\mathbf{V}}\bar{\mathbf{A}}^2\bar{\mathbf{V}}^T, \quad (36)$$

where

$$\bar{\mathbf{A}}^2 = \text{diag} \{ \bar{\lambda}_i^2 \}. \quad (37)$$

$i=1, \dots, n$

Therefore, (36) provides an eigen-decomposition of $\boldsymbol{\Sigma}^2$. By a suitable rearrangement of the columns of $\bar{\mathbf{V}}$ and the diagonal entries of $\bar{\mathbf{A}}$, we can obtain the ordering $\bar{\lambda}_1^2 \geq \dots \geq \bar{\lambda}_n^2 \geq 0$, **without** changing \mathbf{W} . This is easily seen by examining a dyadic expansion of \mathbf{W} in (34), namely

$$\mathbf{W} = \sum_{i=1}^n \bar{\lambda}_i \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^T,$$

where $\bar{\mathbf{v}}_i$ for $i = 1, \dots, n$ are the columns of $\bar{\mathbf{V}}$. Clearly, \mathbf{W} is invariant under order of summation. Since $\boldsymbol{\Sigma}^2$ is diagonal, its diagonal elements are also its eigenvalues. Since the diagonal entries in $\boldsymbol{\Sigma}^2$ are in decreasing order, it must be that $\boldsymbol{\Sigma}^2 = \bar{\mathbf{A}}^2$. Now, the distinct eigenvalues of $\boldsymbol{\Sigma}^2$ are $\sigma_1^2 > \dots > \sigma_m^2 \geq 0$, each with multiplicity n_i for $i = 1, \dots, m$. Correspondingly, $\bar{\mathbf{A}}$ can be written as $\bar{\mathbf{A}} = \text{diag}_{i=1, \dots, m} \{ \bar{\mathbf{A}}_i \}$, where $\bar{\mathbf{A}}_i = \text{diag}_{j=1, \dots, n_i} \{ \bar{\lambda}_{i,j} \}$ for $i = 1, \dots, m$. Consequently, it must be that

$$\sigma_i^2 = \bar{\lambda}_{i,j}^2, \quad i = 1, \dots, m, \quad j = 1, \dots, n_i.$$

Therefore, the eigenvalues of \mathbf{W} must satisfy

$$\bar{\lambda}_{i,j} = p_i^j \sigma_i, \quad p_i^j = \pm 1, \quad i = 1, \dots, m, \quad j = 1, \dots, n_i. \quad (38)$$

Consequently, (34) takes the form given in (31) with $\bar{\mathbf{V}}$ in place of $\boldsymbol{\Delta}$. Finally, we note that $\mathbf{V}_0 = \mathbf{1}$ also provides an eigen-decomposition of $\boldsymbol{\Sigma}^2$, namely

$$\boldsymbol{\Sigma}^2 = \mathbf{V}_0 \boldsymbol{\Sigma}^2 \mathbf{V}_0.$$

Therefore, application of Lemma 5 in the appendix, shows that $\bar{\mathbf{V}} = \boldsymbol{\Delta}$ as in (33). This shows that any real symmetric \mathbf{W} satisfying (30) must take the form given in (31). This concludes the proof.

Let us now substitute (29) into the performance index \hat{J} in (21), to get

$$\hat{J} = \text{tr} [\mathbf{SWS}^T] = \text{tr} [\mathbf{S}^T \mathbf{S} \mathbf{W}] = \text{tr} [\mathbf{W}]. \quad (39)$$

Substituting (31) into (39), we find that

$$\hat{J} = \sum_{i=1}^m \sum_{j=1}^{n_i} p_i^j \sigma_i. \quad (40)$$

It is readily seen from (40) that the performance index is maximized by having as many of the p_i^j coefficients as possible positive, while still satisfying the determinant constraint (namely that $\det \mathbf{S} = \det \mathbf{U} \det \mathbf{V}$). If any need to be negative, they should correspond to the smallest eigenvalues. There are now several cases to consider:

1. $\det \mathbf{B} \neq 0$
 - (a) $\det \mathbf{B} > 0$
 - (b) $\det \mathbf{B} < 0$
 - i. $n_m = 1$
 - ii. $n_m > 1$
2. $\det \mathbf{B} = 0$
 - (a) $\text{rank}[\mathbf{B}] = n - 1$
 - (b) $\text{rank}[\mathbf{B}] < n - 1$

Each of these cases will be investigated separately in the following subsections, and all findings are summarized in Theorem 1.

Case 1, $\det \mathbf{B} \neq 0$

In this case, from (17) we find that Σ has full rank, and therefore $\sigma_m > 0$. Furthermore, since $\det \Sigma > 0$,

$$\text{sign}[\det \mathbf{B}] = \det \mathbf{U} \det \mathbf{V}. \quad (41)$$

Consequently, from (31) it follows that any \mathbf{W} satisfying (30) is non-singular. Therefore, from (29) and Lemma 1 it follows that any possible maximizing solution of (21) must have the form

$$\begin{aligned} \mathbf{S} &= \mathbf{W}^{-1} \Sigma = \text{diag}_{i=1, \dots, m} \{ \Delta_i \bar{\Sigma}_i^{-1} \Delta_i^T (\sigma_i \mathbf{1}_{n_i \times n_i}) \}, \\ &= \Delta \text{diag}_{i=1, \dots, m} \{ \mathbf{P}_i \} \Delta^T, \end{aligned} \quad (42)$$

where

$$\mathbf{P}_i = \text{diag}_{j=1, \dots, n_i} \{ p_i^j \}, \quad p_i^j = \pm 1, \quad j = 1, \dots, n_i, \quad (43)$$

for $i = 1, \dots, m$, and Δ satisfies (33). Using (42), it is readily found by direct calculation that $\mathbf{S} \mathbf{S}^T = \mathbf{1}$. It now remains to identify the coefficients p_i^j in (43)

that maximize \hat{J} in (40), while ensuring that $\det \mathbf{S} = \det \mathbf{U} \det \mathbf{V}$. From (42), it is seen that this leads to the requirement

$$\prod_{i=1}^m \prod_{j=1}^{n_i} p_i^j = \det \mathbf{U} \det \mathbf{V}. \quad (44)$$

Case 1 (a), $\det \mathbf{B} > 0$

From (41) and (44), it follows that $\prod_{i=1}^m \prod_{j=1}^{n_i} p_i^j = 1$ (since $\det \mathbf{\Sigma} > 0$). From (40), it is now seen that \hat{J} is uniquely maximized by the choice $p_i^j = 1$ for $j = 1, \dots, n_i$, $i = 1, \dots, m$. Therefore, from (42), it follows that the unique global maximizing solution to (21) is

$$\mathbf{S} = \mathbf{1}. \quad (45)$$

From (20), it follows that the unique global minimizing solution to Problem 1 is

$$\mathbf{C} = \mathbf{V}\mathbf{U}^T. \quad (46)$$

This solution is the same as the one presented in [9] (which only considered the case $\det \mathbf{B} > 0$).

Case 1 (b), $\det \mathbf{B} < 0$

From (41) and (44), it follows that $\prod_{i=1}^m \prod_{j=1}^{n_i} p_i^j = -1$, which means that at least one of the p_i^j should be negative. Since $\sigma_1 > \dots > \sigma_m > 0$, we readily find from (40) that there are n_m choices for p_i^j (the multiplicity of σ_m) leading to maximizing solutions of the performance index \hat{J} . These are given by $p_i^j = 1$ for $j = 1, \dots, n_i$, $i = 1, \dots, m-1$, and

$$p_{m,k}^j = \begin{cases} -1, & j = k, \\ 1, & j \in \{1, \dots, n_m\} - \{k\}, \end{cases} \quad (47)$$

for $k = 1, \dots, n_m$. Let us now examine the corresponding \mathbf{S} , given in (42). We have

$$\mathbf{S} = \text{diag} \left\{ \text{diag}_{i=1, \dots, m-1} \{ \mathbf{1}_{n_i \times n_i} \}, \mathbf{\Delta}_m \bar{\mathbf{P}}_m^k \mathbf{\Delta}_m^T \right\}, \quad (48)$$

where

$$\bar{\mathbf{P}}_m^k = \text{diag}_{j=1, \dots, n_m} \{ p_{m,k}^j \},$$

and $\mathbf{\Delta}_m \in \mathbb{R}^{n_m \times n_m}$ is a real orthonormal matrix. Now, by a suitable re-ordering of the columns of $\mathbf{\Delta}_m$ and the diagonal entries of $\bar{\mathbf{P}}_m^k$, we can find $\bar{\mathbf{\Delta}}_m$ (also real and orthonormal), such that

$$\mathbf{\Delta}_m \bar{\mathbf{P}}_m^k \mathbf{\Delta}_m^T = \bar{\mathbf{\Delta}}_m \bar{\mathbf{P}}_m^{n_m} \bar{\mathbf{\Delta}}_m^T.$$

Note that $\bar{\mathbf{P}}_m^{n_m} = \text{diag}\{1, \dots, 1, -1\}$. Therefore, since $\mathbf{\Delta}_m$ is a free-parameter in (48), we do not need to distinguish between the n_m separate possibilities for $p_{m,k}^j$ given in (47).

Consequently, all maximizing solutions for \mathbf{S} may be generated according to

$$\mathbf{S} = \text{diag} \left\{ \text{diag}_{i=1, \dots, m-1} \{ \mathbf{1}_{n_i \times n_i} \}, \mathbf{\Delta}_m \hat{\mathbf{P}}_m \mathbf{\Delta}_m^T \right\}, \quad (49)$$

where $\mathbf{\Delta}_m$ satisfies $\mathbf{\Delta}_m^T \mathbf{\Delta}_m = \mathbf{1}_{n_m \times n_m}$, and $\hat{\mathbf{P}}_m = \text{diag}\{1, \dots, 1, -1\}$. There are now two sub-cases to consider:

Case 1 (b) i: $n_m = 1$ (the minimum singular value (σ_m) of \mathbf{B} is distinct)

In this case, $\mathbf{\Delta}_m = \pm 1$, $\hat{\mathbf{P}}_m = -1$, and the solution in (49) becomes

$$\mathbf{S} = \text{diag}\{\mathbf{1}_{(n-1) \times (n-1)}, -1\}, \quad (50)$$

which is unique. Consequently, from (20), the minimizing solution \mathbf{C} for Problem 1 is unique also, and is given by

$$\mathbf{C} = \mathbf{V} \text{diag}\{\mathbf{1}_{(n-1) \times (n-1)}, -1\} \mathbf{U}^T. \quad (51)$$

Case 1 (b) ii: $n_m > 1$ (the minimum singular value (σ_m) of \mathbf{B} is repeated)

In this case, \mathbf{S} in (49) is non-unique, and consequently, the solutions of Problem 1 are non-unique also. From (20) and (49), the set of all minimizing solutions of Problem 1 are given by

$$\mathbf{C} = \mathbf{V} \text{diag} \left\{ \text{diag}_{i=1, \dots, m-1} \{ \mathbf{1}_{n_i \times n_i} \}, \mathbf{\Delta}_m \hat{\mathbf{P}}_m \mathbf{\Delta}_m^T \right\} \mathbf{U}^T, \quad (52)$$

where $\mathbf{\Delta}_m$ satisfies $\mathbf{\Delta}_m^T \mathbf{\Delta}_m = \mathbf{1}_{n_m \times n_m}$. By application of Lemma 6 (see the appendix), it is readily found that the set of all minimizing \mathbf{C} in (52) is equivalently given by

$$\mathbf{C} = \mathbf{V} \text{diag}\{\mathbf{1}_{(n-1) \times (n-1)}, -1\} \mathbf{U}^T, \quad (53)$$

for all real orthonormal \mathbf{V} and \mathbf{U} satisfying

$$\mathbf{B} = \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T.$$

Remark From (41), the global minimizing solutions to Problem 1 in (46), (51) and (53) for all cases of $\det \mathbf{B} \neq 0$ can all be written as

$$\mathbf{C} = \mathbf{V} \text{diag}\{\mathbf{1}_{(n-1) \times (n-1)}, \det \mathbf{V} \det \mathbf{U}\} \mathbf{U}^T, \quad (54)$$

where \mathbf{V} and \mathbf{U} satisfy the singular value decomposition $\mathbf{B} = \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T$. It will subsequently be shown that when $\det \mathbf{B} = 0$, the solutions to Problem 1 can also be written in the form of (54).

Case 2: $\det \mathbf{B} = 0$

In this case, \mathbf{B} does not have full-rank, which implies that $\sigma_m = 0$. The matrix \mathbf{W} in Lemma 1 and the matrix $\mathbf{\Sigma}$ now take the form

$$\mathbf{W} = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{\Sigma} = \begin{bmatrix} \bar{\mathbf{A}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad (55)$$

where

$$\mathbf{A} = \text{diag}_{i=1, \dots, m-1} \{ \mathbf{\Delta}_i \bar{\mathbf{\Sigma}}_i \mathbf{\Delta}_i^T \}, \quad \bar{\mathbf{A}} = \text{diag}_{i=1, \dots, m-1} \{ \sigma_i \mathbf{1}_{n_i \times n_i} \}, \quad (56)$$

the matrices $\bar{\mathbf{\Sigma}}_i$ are given by (32), and $\mathbf{\Delta}_i \in \mathbb{R}^{n_i \times n_i}$ are real orthonormal matrices for $i = 1, \dots, m-1$. Let us correspondingly partition \mathbf{S} as

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}_{11} & \mathbf{S}_{12} \\ \mathbf{S}_{21} & \mathbf{S}_{22} \end{bmatrix}, \quad (57)$$

where $\mathbf{S}_{11} \in \mathbb{R}^{(n-n_m) \times (n-n_m)}$, $\mathbf{S}_{12} \in \mathbb{R}^{(n-n_m) \times n_m}$, $\mathbf{S}_{21} \in \mathbb{R}^{n_m \times (n-n_m)}$ and $\mathbf{S}_{22} \in \mathbb{R}^{n_m \times n_m}$. Using (55), the necessary condition in (29) leads to

$$\mathbf{A} \mathbf{S}_{11} = \bar{\mathbf{A}}, \quad \mathbf{A} \mathbf{S}_{12} = \mathbf{0},$$

and using (56), this gives

$$\mathbf{S}_{11} = \text{diag}_{i=1, \dots, m-1} \{ \mathbf{\Delta}_i \mathbf{P}_i \mathbf{\Delta}_i^T \}, \quad \mathbf{S}_{12} = \mathbf{0}, \quad (58)$$

where \mathbf{P}_i are given in (43) and \mathbf{S}_{21} and \mathbf{S}_{22} remain free parameters. Since \mathbf{S} must be orthonormal, we have

$$\mathbf{S}^T \mathbf{S} = \begin{bmatrix} \mathbf{S}_{11}^T & \mathbf{S}_{21}^T \\ \mathbf{0} & \mathbf{S}_{22}^T \end{bmatrix} \begin{bmatrix} \mathbf{S}_{11} & \mathbf{0} \\ \mathbf{S}_{21} & \mathbf{S}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix},$$

which leads to

$$\mathbf{S}_{11}^T \mathbf{S}_{11} + \mathbf{S}_{21}^T \mathbf{S}_{21} = \mathbf{1}, \quad \mathbf{S}_{22}^T \mathbf{S}_{22} = \mathbf{1}, \quad \mathbf{S}_{21}^T \mathbf{S}_{22} = \mathbf{0}.$$

From (58), we have $\mathbf{S}_{11}^T \mathbf{S}_{11} = \mathbf{1}$, leading to $\mathbf{S}_{21}^T \mathbf{S}_{21} = \mathbf{0}$, and consequently $\mathbf{S}_{21} = \mathbf{0}$. Using (58), the determinant of \mathbf{S} becomes

$$\det \mathbf{S} = \det \mathbf{S}_{11} \det \mathbf{S}_{22} = \left(\prod_{i=1}^{m-1} \prod_{j=1}^{n_i} p_i^j \right) \det \mathbf{S}_{22}.$$

Since $p_i^j = \pm 1$ for $j = 1, \dots, n_i$, $i = 1, \dots, m-1$, the determinant requirement $\det \mathbf{S} = \det \mathbf{V} \det \mathbf{U}$ therefore is

$$\det \mathbf{S}_{22} = \det \mathbf{V} \det \mathbf{U} \left(\prod_{i=1}^{m-1} \prod_{j=1}^{n_i} p_i^j \right). \quad (59)$$

With $\sigma_m = 0$, the performance index in (40) becomes

$$\hat{J} = \sum_{i=1}^{m-1} \sum_{j=1}^{n_i} p_i^j \sigma_i, \quad (60)$$

which is independent of \mathbf{S}_{22} . The performance index \hat{J} is therefore globally maximized by selecting $p_i^j = 1$ for $i = 1, \dots, m-1$, $j = 1, \dots, n_i$. From (58), this leads to

$$\mathbf{S}_{11} = \text{diag}_{i=1, \dots, m-1} \{ \mathbf{1}_{n_i \times n_i} \},$$

while the matrix \mathbf{S}_{22} is then a free parameter that can be used to enforce the determinant constraint. Therefore, we conclude that the set of all global maximizing solutions \mathbf{S} may be generated according to

$$\mathbf{S} = \text{diag} \left\{ \text{diag}_{i=1, \dots, m-1} \{ \mathbf{1}_{n_i \times n_i} \}, \mathbf{S}_{22} \right\}, \quad (61)$$

where $\mathbf{S}_{22} \in \mathbb{R}^{n_m \times n_m}$ is orthonormal and has determinant $\det \mathbf{S}_{22} = \det \mathbf{V} \det \mathbf{U}$.

Case 2 (a): $\text{rank}[\mathbf{B}] = n - 1$ ($n_m = 1$)

In this case, \mathbf{S}_{22} in (61) becomes a scalar, and the determinant condition in (61) becomes

$$S_{22} = \det \mathbf{V} \det \mathbf{U}.$$

Note that since $\det \mathbf{V} \det \mathbf{U} = \pm 1$, the orthonormality condition in (61) is automatically satisfied. Therefore, it follows from (61) that the unique global maximizing solution of (21) is

$$\mathbf{S} = \text{diag} \{ \mathbf{1}_{(n-1) \times (n-1)}, \det \mathbf{V} \det \mathbf{U} \}. \quad (62)$$

Consequently, from (20), the minimizing solution \mathbf{C} for Problem 1 is unique also, and is given by

$$\mathbf{C} = \mathbf{V} \text{diag} \{ \mathbf{1}_{(n-1) \times (n-1)}, \det \mathbf{V} \det \mathbf{U} \} \mathbf{U}^T. \quad (63)$$

We note once again that (63) is identical in form to (54).

Case 2 (b): $\text{rank}[\mathbf{B}] < n - 1$ ($n_m > 1$)

In this case, since $n_m > 1$, the set of all orthonormal $\mathbf{S}_{22} \in \mathbb{R}^{n_m \times n_m}$ satisfying the determinant requirement in (61) is infinite, and hence the set of global maximizing solutions for (21) is non-unique. Consequently, from (20) and (61), the set of all global solutions of Problem 1 is non-unique and is given by

$$\mathbf{C} = \mathbf{V} \text{diag} \left\{ \text{diag}_{i=1, \dots, m-1} \{ \mathbf{1}_{n_i \times n_i} \}, \mathbf{S}_{22} \right\} \mathbf{U}^T, \quad (64)$$

where $\mathbf{S}_{22} \in \mathbb{R}^{n_m \times n_m}$ is orthonormal with determinant $\det \mathbf{S}_{22} = \det \mathbf{V} \det \mathbf{U}$.

We shall now show that the set of solutions in (64) may once again be put into the form of (54). To this end, define the matrix

$$\hat{\mathbf{A}} = \text{diag}\{\mathbf{1}_{(n_m-1) \times (n_m-1)}, \det \mathbf{V} \det \mathbf{U}\}. \quad (65)$$

Clearly, $\hat{\mathbf{A}}$ is orthonormal and $\det \hat{\mathbf{A}} = \det \mathbf{V} \det \mathbf{U}$. As such, setting $\mathbf{S}_{22} = \hat{\mathbf{A}}$ in (64) yields a minimizing solution of Problem 1. With this choice, (64) takes the form of (54). Since the choice of \mathbf{V} and \mathbf{U} satisfying (17) was arbitrary, it follows that (54) leads to a minimizing solution of Problem 1 for any choice of \mathbf{V} and \mathbf{U} satisfying (17). Conversely, for a given \mathbf{S}_{22} in (64), define the orthonormal matrices

$$\bar{\mathbf{V}} = \mathbf{V} \text{diag} \left\{ \text{diag}_{i=1, \dots, m-1} \{ \mathbf{1}_{n_i \times n_i} \}, \mathbf{S}_{22} \right\}, \quad \bar{\mathbf{U}} = \mathbf{U} \text{diag} \left\{ \text{diag}_{i=1, \dots, m-1} \{ \mathbf{1}_{n_i \times n_i} \}, \hat{\mathbf{A}} \right\},$$

such that

$$\mathbf{C} = \bar{\mathbf{V}} \text{diag} \left\{ \text{diag}_{i=1, \dots, m-1} \{ \mathbf{1}_{n_i \times n_i} \}, \hat{\mathbf{A}} \right\} \bar{\mathbf{U}}^T. \quad (66)$$

By Lemma 6 in the appendix, $\bar{\mathbf{V}}$ and $\bar{\mathbf{U}}$ satisfy $\mathbf{B} = \bar{\mathbf{V}} \boldsymbol{\Sigma} \bar{\mathbf{U}}^T$. Now,

$$\begin{aligned} \det \bar{\mathbf{V}} \det \bar{\mathbf{U}} &= \det \mathbf{V} \det \mathbf{U} \det \mathbf{S}_{22} \det \hat{\mathbf{A}}, \\ &= (\det \mathbf{V} \det \mathbf{U})^3, \\ &= \det \mathbf{V} \det \mathbf{U}. \end{aligned} \quad (67)$$

Therefore, (65), (66) and (67) show that for a given \mathbf{S}_{22} in (64), one may equivalently write \mathbf{C} in (64) as

$$\mathbf{C} = \bar{\mathbf{V}} \text{diag} \{ \mathbf{1}_{(n-1) \times (n-1)}, \det \bar{\mathbf{V}} \det \bar{\mathbf{U}} \} \bar{\mathbf{U}}^T, \quad (68)$$

for some orthonormal $\bar{\mathbf{V}}$ and $\bar{\mathbf{U}}$ satisfying $\mathbf{B} = \bar{\mathbf{V}} \boldsymbol{\Sigma} \bar{\mathbf{U}}^T$, which also has the form of (54).

The findings are now summarized in the following Theorem.

Theorem 1 *The set of all solutions to Problem 1 are generated by*

$$\mathbf{C} = \mathbf{V} \text{diag} \{ \mathbf{1}_{(n-1) \times (n-1)}, \det \mathbf{V} \det \mathbf{U} \} \mathbf{U}^T, \quad (69)$$

where \mathbf{V} and \mathbf{U} satisfy the singular value decomposition of \mathbf{B} such that

$$\mathbf{B} = \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T,$$

where $\mathbf{V}^T \mathbf{V} = \mathbf{1}$, $\mathbf{U}^T \mathbf{U} = \mathbf{1}$ and $\boldsymbol{\Sigma} = \text{diag}_{i=1, \dots, n} \{ \sigma_i \}$ with $\sigma_1 \geq \dots \geq \sigma_n \geq 0$.

Furthermore, the solution is unique if

1. $\det \mathbf{B} > 0$,
2. $\det \mathbf{B} < 0$ and the minimum singular value of \mathbf{B} is distinct,
3. $\text{rank}[\mathbf{B}] = n - 1$.

Note that (69) was derived in [10] and [12] by different means for the special case of $n = 3$ ($\mathbf{C} \in SO(3)$). Reference [9] obtains this result in the case $\det \mathbf{B} > 0$ for $SO(n)$. Reference [11] obtains this Theorem for $SO(n)$, but erroneously concludes uniqueness of solutions for $\det \mathbf{B} < 0$, without requiring the minimum singular value of \mathbf{B} to be distinct. Finally, we have characterized the entire set of solutions of Problem 1, which was not done previously.

Suppose now that we modify Problem 1 by reversing the determinant requirement.

Problem 4 *Given N vectors, $\mathbf{s}_{b,k} \in \mathbb{R}^n$, with corresponding vectors $\mathbf{s}_{I,k} \in \mathbb{R}^n$, find the matrix $\mathbf{C} \in \{\mathbf{C} \in \mathbb{R}^{n \times n} : \mathbf{C}^T \mathbf{C} = \mathbf{1}, \det \mathbf{C} = -1\}$, to minimize the cost function*

$$J = \sum_{k=1}^N w_k (\mathbf{s}_{b,k} - \mathbf{C} \mathbf{s}_{I,k})^T (\mathbf{s}_{b,k} - \mathbf{C} \mathbf{s}_{I,k}), \quad (70)$$

where $0 < w_k < \infty$ are positive weights for $k = 1, \dots, N$.

An example of where this would occur is when $\mathbf{C} \in \mathbb{R}^{3 \times 3}$ is a rotation matrix describing the transformation between right and left-handed reference frames, or vice-versa. Repeating the analysis for Problem 1, we readily obtain the result

Theorem 2 *The set of all solutions to Problem 4 are generated by*

$$\mathbf{C} = \mathbf{V} \text{diag} \{ \mathbf{1}_{(n-1) \times (n-1)}, -\det \mathbf{V} \det \mathbf{U} \} \mathbf{U}^T, \quad (71)$$

where \mathbf{V} and \mathbf{U} satisfy the singular value decomposition of \mathbf{B} such that

$$\mathbf{B} = \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T,$$

where $\mathbf{V}^T \mathbf{V} = \mathbf{1}$, $\mathbf{U}^T \mathbf{U} = \mathbf{1}$ and $\boldsymbol{\Sigma} = \text{diag}_{i=1, \dots, n} \{ \sigma_i \}$ with $\sigma_1 \geq \dots \geq \sigma_n \geq 0$.

Furthermore, the solution is unique if

1. $\det \mathbf{B} < 0$,
2. $\det \mathbf{B} > 0$ and the minimum singular value of \mathbf{B} is distinct,
3. $\text{rank}[\mathbf{B}] = n - 1$.

Finally, suppose we remove the determinant requirement altogether to obtain

Problem 5 *Given N vectors, $\mathbf{s}_{b,k} \in \mathbb{R}^n$, with corresponding vectors $\mathbf{s}_{I,k} \in \mathbb{R}^n$, find the matrix $\mathbf{C} \in O(n)$, to minimize the cost function*

$$J = \sum_{k=1}^N w_k (\mathbf{s}_{b,k} - \mathbf{C} \mathbf{s}_{I,k})^T (\mathbf{s}_{b,k} - \mathbf{C} \mathbf{s}_{I,k}), \quad (72)$$

where $0 < w_k < \infty$ are positive weights for $k = 1, \dots, N$.

Then, repeating the analysis, we readily obtain the following Theorem, which has also been presented in [7].

Theorem 3 *The set of all solutions to Problem 5 are generated by*

$$\mathbf{C} = \mathbf{V}\mathbf{U}^T, \quad (73)$$

where \mathbf{V} and \mathbf{U} satisfy the singular value decomposition of \mathbf{B} such that

$$\mathbf{B} = \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T,$$

where $\mathbf{V}^T\mathbf{V} = \mathbf{1}$, $\mathbf{U}^T\mathbf{U} = \mathbf{1}$ and $\mathbf{\Sigma} = \text{diag}_{i=1,\dots,n} \{\sigma_i\}$ with $\sigma_1 \geq \dots \geq \sigma_n \geq 0$.

Furthermore, the solution is unique if and only if $\text{rank}[\mathbf{B}] = n$.

Note that when $\text{rank}[\mathbf{B}] = n - 1$, there are two solutions to Problem 5, and when $\text{rank}[\mathbf{B}] < n - 1$, there are infinite solutions to Problem 5. Comparing Theorems 1 and 2 to Theorem 3, we see that it is the determinant conditions in Problems 1 and 4 that guarantee uniqueness when $\text{rank}[\mathbf{B}] = n - 1$.

Finally, (41) shows that the minimizing solution \mathbf{C} of Problem 5 when $\text{rank}[\mathbf{B}] = n$ has

$$\det \mathbf{C} = \text{sign}[\det \mathbf{B}].$$

As such, comparing (69) to (73), we see that when $\det \mathbf{B} > 0$, the solution of Problem 5 yields the solution of Problem 1. Likewise, comparing (71) to (73), we see that when $\det \mathbf{B} < 0$, the solution of Problem 5 yields the solution of Problem 4.

Example 1 *As an interesting example, suppose that we are given an orthogonal matrix $\mathbf{D} \in \mathbb{R}^{2 \times 2}$ with determinant $\det \mathbf{D} = -1$, that is, $\mathbf{D} \in O(2) - SO(2)$. Suppose that it is desired to find the matrix $\mathbf{C} \in SO(2)$ to solve Problem 2. First, we note that \mathbf{D} has singular values $\sigma_1 = \sigma_2 = 1$, and that $\mathbf{V} = \mathbf{D}$, $\mathbf{\Sigma} = \mathbf{1}_{2 \times 2}$, and $\mathbf{U} = \mathbf{1}_{2 \times 2}$ satisfy a singular value decomposition of $\mathbf{D} = \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T$. Consequently, from (52), it follows that the set of all minimizing $\mathbf{C} \in SO(2)$ is given by*

$$\mathbf{C} = \mathbf{D}\mathbf{\Delta}\mathbf{P}\mathbf{\Delta}^T, \quad (74)$$

where $\mathbf{\Delta} \in \mathbb{R}^{2 \times 2}$ is an orthonormal matrix, and $\mathbf{P} = \text{diag}\{1, -1\}$.

Next, we note that all matrices $\mathbf{D} \in O(2) - SO(2)$ take the form of

$$\mathbf{D} = \begin{bmatrix} \cos \theta & -\sin \theta \\ -\sin \theta & -\cos \theta \end{bmatrix}, \quad (75)$$

for some angle $\theta \in \mathbb{R}$, while all matrices $\mathbf{C} \in SO(2)$ take the form

$$\mathbf{C} = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}, \quad (76)$$

for some angle $\theta \in \mathbb{R}$. Now, given any $\mathbf{C} \in SO(2)$, the equation

$$\mathbf{C} = \mathbf{D}\mathbf{S}, \quad (77)$$

is solvable for $\mathbf{S} \in O(2) - SO(2)$. We shall now show that for any $\mathbf{S} \in O(2) - SO(2)$, the equation

$$\mathbf{S} = \mathbf{\Delta}\mathbf{P}\mathbf{\Delta}^T, \quad (78)$$

is solvable for $\Delta \in O(2)$. We know from (75) that \mathbf{S} must take the form

$$\mathbf{S} = \begin{bmatrix} \cos \theta & -\sin \theta \\ -\sin \theta & -\cos \theta \end{bmatrix}, \quad (79)$$

for some angle $\theta \in \mathbb{R}$. Setting $\Delta \in O(2)$ equal to

$$\Delta = \begin{bmatrix} \cos \theta/2 & \sin \theta/2 \\ -\sin \theta/2 & \cos \theta/2 \end{bmatrix}, \quad (80)$$

it is readily found using double angle trigonometric identities that Δ does indeed solve (79). Consequently, given any $\mathbf{C} \in SO(2)$, it is possible to find a $\Delta \in O(2)$ satisfying (74). Therefore, the set of solutions to Problem 2 in this case is in fact all of $SO(2)$.

Matrix Square-Root Solution when $\det \mathbf{B} \neq 0$

Another useful form of the solution to Problem 1 is the matrix square-root solution, which was originally presented in [6] for problems in $SO(n)$, and in [8] for problems in $SO(3)$. However, conditions for uniqueness of solutions were not comprehensively studied in the aforementioned reference. Using the singular-value decomposition solution from the previous section, it is straight forward to obtain the the complete set of solutions in terms of the matrix square-root. Note that the matrix square-root solution is only defined when $\det \mathbf{B} \neq 0$.

To this end, given a singular value decomposition of \mathbf{B} in (17), define the matrices

$$\mathbf{W}_l = \mathbf{V} \text{diag} \left\{ \text{diag}_{i=1, \dots, m-1} \{ \sqrt{\lambda_i} \mathbf{1}_{n_i \times n_i} \}, \sqrt{\lambda_m} \hat{\mathbf{A}}_m \right\} \mathbf{V}^T, \quad (81)$$

and

$$\mathbf{W}_r = \mathbf{U} \text{diag} \left\{ \text{diag}_{i=1, \dots, m-1} \{ \sqrt{\lambda_i} \mathbf{1}_{n_i \times n_i} \}, \sqrt{\lambda_m} \hat{\mathbf{A}}_m \right\} \mathbf{U}^T, \quad (82)$$

where

$$\hat{\mathbf{A}}_m = \text{diag} \{ \mathbf{1}_{(n_m-1) \times (n_m-1)}, \text{sign}[\det \mathbf{B}] \},$$

and $\lambda_i = \sigma_i^2$ for $i = 1, \dots, m$. Using (81), (82) and (17), it is easy to see that \mathbf{W}_l and \mathbf{W}_r are real symmetric square-root matrices of $\mathbf{B}\mathbf{B}^T$ and $\mathbf{B}^T\mathbf{B}$, respectively. Furthermore, λ_i are the distinct eigenvalues of $\mathbf{B}\mathbf{B}^T$ and $\mathbf{B}^T\mathbf{B}$, while the columns of \mathbf{V} are their corresponding eigenvectors for $\mathbf{B}\mathbf{B}^T$ and the columns of \mathbf{U} are their corresponding eigenvectors for $\mathbf{B}^T\mathbf{B}$.

Using (17), (81) and (82), together with (41) and (69), it is straightforward to verify that

$$\mathbf{C} = \mathbf{W}_l^{-1} \mathbf{B} = \mathbf{B} \mathbf{W}_r^{-1} \quad (83)$$

Finally, combining Lemmas 5 and 6 in the appendix, it is straightforward to show that $\mathbf{V} \in \mathbb{R}^{n \times n}$ is a real orthonormal matrix diagonalizing $\mathbf{B}\mathbf{B}^T =$

$\mathbf{V}\boldsymbol{\Sigma}^2\mathbf{V}^T$ if and only if there exists a real orthonormal matrix $\mathbf{U} \in \mathbb{R}^n$ such that $\mathbf{B} = \mathbf{V}\boldsymbol{\Sigma}\mathbf{U}^T$. Likewise, $\mathbf{U} \in \mathbb{R}^{n \times n}$ is a real orthonormal matrix diagonalizing $\mathbf{B}^T\mathbf{B} = \mathbf{U}\boldsymbol{\Sigma}^2\mathbf{U}^T$ if and only if there exists a real orthonormal matrix $\mathbf{V} \in \mathbb{R}^n$ such that $\mathbf{B} = \mathbf{V}\boldsymbol{\Sigma}\mathbf{U}^T$. Therefore, it is concluded that the entire set of global minimizing solutions for Problem 1 is generated by (83) together with either (81) or (82), where $\mathbf{V} \in \mathbb{R}^{n \times n}$ and $\mathbf{U} \in \mathbb{R}^{n \times n}$ are real orthonormal matrices diagonalizing $\mathbf{B}\mathbf{B}^T$ or $\mathbf{B}^T\mathbf{B}$ respectively, where their distinct eigenvalues are ordered according to $\lambda_1 > \dots > \lambda_m > 0$, each with multiplicity n_i .

Remark *The matrix square-root solutions in (83) are often written as*

$$\mathbf{C} = (\mathbf{B}\mathbf{B}^T)^{-1/2}\mathbf{B} = \mathbf{B}(\mathbf{B}^T\mathbf{B})^{-1/2}.$$

However, this notation is ambiguous as to which matrix square-root is used. The correct choices of matrix square root are provided in (81) and (82).

An LMI Method of Solution when $\det \mathbf{B} > 0$

Reference [16] (by the authors of this paper) presents an LMI method of solution to Wahba's problem on $SO(3)$, when $\det \mathbf{B} > 0$. This section presents the general case on $SO(n)$, which is a simple generalization of the results presented in [16].

First, we note that Problems 1, 4 and 5 may all be reduced to a maximization problem of the form given in (22) (with $(\mathbf{C}, \mathbf{B}^T)$ in place of $(\mathbf{S}, \boldsymbol{\Sigma})$), with the addition of a determinant constraint in the cases of Problems 1 and 4 (see for example (6), for Problem 1). The constraint set in (22) is non-convex, and in this section we show how it can be relaxed to a convex constraint set.

Consider the problem

Problem 6 *Let $\mathbf{B} \in \mathbb{R}^{n \times n}$.*

$$\text{maximize } \hat{J} = \text{tr}[\mathbf{C}\mathbf{B}^T] \text{ subject to } \|\mathbf{C}\| \leq 1, \quad (84)$$

where $\|\mathbf{A}\| = \sigma_{\max}(\mathbf{A})$ is the induced 2-norm of the matrix \mathbf{A} , and σ_{\max} the maximum singular value.

This problem appears similar to the maximization problem in (22) (with $(\mathbf{C}, \mathbf{B}^T)$ in place of $(\mathbf{S}, \boldsymbol{\Sigma})$), except that the constraint set is relaxed. We note that $\|\mathbf{C}\| = 1$ for all orthonormal \mathbf{C} , so the constraint set in Problem 6 contains the constraint set in (22).

Next, we note that the constraint set is convex. Indeed, consider any \mathbf{C}_1 and \mathbf{C}_2 satisfying $\|\mathbf{C}_1\| \leq 1$ and $\|\mathbf{C}_2\| \leq 1$, and form the convex combination $\alpha\mathbf{C}_1 + (1 - \alpha)\mathbf{C}_2$ for some $\alpha \in [0, 1]$. Then, using the properties of the matrix norm,

$$\|\alpha\mathbf{C}_1 + (1 - \alpha)\mathbf{C}_2\| \leq \alpha\|\mathbf{C}_1\| + (1 - \alpha)\|\mathbf{C}_2\| \leq 1.$$

As a consequence, since the performance index in (84) is linear (and hence convex) any local maximizing solution of (84) must be a global maximizing solution.

As before, the constraint set in Problem 6 is compact, so a global minimizing solution exists. Furthermore, since the performance index \hat{J} is linear in \mathbf{C} , the maximizing solution must lie on the boundary of the constraint set. That is, the maximizing solution must satisfy $\|\mathbf{C}\| = 1$. Note however, that this by itself does not guarantee that the maximizing solution is orthonormal. For example, $\mathbf{C} = \text{diag}\{\mathbf{0}_{(n-1) \times (n-1)}, 1\}$ has unit norm, but is not orthonormal. It will be demonstrated in the following that when \mathbf{B} has full rank, then the maximizing solution of Problem 6 is indeed orthonormal, and solves (22) (with $(\mathbf{C}, \mathbf{B}^T)$ in place of $(\mathbf{S}, \boldsymbol{\Sigma})$). Consequently, when \mathbf{B} has full rank, the solution to Problem 6 also solves Problem 5, and depending on the sign of $\det \mathbf{B}$, solves one of Problems 1 and 4. In this manner, it can be shown that the non-convex maximization problem in (22) (with $(\mathbf{C}, \mathbf{B}^T)$ in place of $(\mathbf{S}, \boldsymbol{\Sigma})$) can be replaced by a convex one, by relaxing the constraint. This may have applications in other optimization problems with orthonormality constraints.

As before, let us consider a singular value decomposition of \mathbf{B} , given by

$$\mathbf{B} = \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T, \quad (85)$$

where $\mathbf{V}^T \mathbf{V} = \mathbf{1}$, $\mathbf{U}^T \mathbf{U} = \mathbf{1}$ and $\boldsymbol{\Sigma} = \text{diag}_{i=1, \dots, n} \{\sigma_i\}$ with $\sigma_1 \geq \dots \geq \sigma_n \geq 0$. As before, consider the change of variables in (18). This is clearly a bijective mapping of the set $\{\mathbf{C} \in \mathbb{R}^{n \times n} : \|\mathbf{C}\| \leq 1\}$ onto itself. As before, the inverse mapping is given by (20).

Then, Problem 6 is equivalent to (compare to (22))

$$\text{maximize } \hat{J} = \text{tr}[\mathbf{S} \boldsymbol{\Sigma}] \text{ subject to } \|\mathbf{S}\| \leq 1. \quad (86)$$

Let us denote the ij^{th} term of \mathbf{S} by S_{ij} . Then, the performance index in (86) becomes

$$\hat{J} = \sum_{i=1}^n S_{ii} \sigma_i. \quad (87)$$

By Corollary 1 (see the appendix), if $\|\mathbf{S}\| \leq 1$ we must have $|S_{ii}| \leq 1$. As such, we have

$$\hat{J}(\mathbf{S}) \leq \sum_{i=1}^n \sigma_i, \quad \forall \mathbf{S} \text{ satisfying } \|\mathbf{S}\| \leq 1. \quad (88)$$

Noting that $\mathbf{S} = \mathbf{1}$ is a member of the constraint set, the upper bound in (88) is in fact the global minimum for (86). Therefore, all minimizing \mathbf{S} of (86) must have

$$S_{ii} = 1 \text{ if } \sigma_i > 0.$$

Correspondingly, by Proposition 1 (see the appendix) they must have

$$S_{ij} = S_{ji} = 0 \text{ for } j \neq i \text{ if } \sigma_i > 0.$$

Hence, if $\sigma_i > 0$ for $i = 1, \dots, m$ and $\sigma_i = 0$ for $i = m + 1, \dots, n$, any maximizing \mathbf{S} of (86) takes the form

$$\mathbf{S} = \begin{bmatrix} \mathbf{1}_{m \times m} & \mathbf{0}_{m \times (n-m)} \\ \mathbf{0}_{(n-m) \times m} & \mathbf{S}_{22} \end{bmatrix}, \quad (89)$$

where $\mathbf{S}_{22} \in \mathbb{R}^{(n-m) \times (n-m)}$. By the norm constraint $\|\mathbf{S}\| \leq 1$, it must satisfy $\|\mathbf{S}_{22}\| \leq 1$. Clearly (from (87)), \mathbf{S}_{22} has no effect on the performance index, and hence may arbitrarily be chosen subject to the norm-constraint.

Case 1: $\text{rank}[\mathbf{B}] = n$

When \mathbf{B} has full rank, $\sigma_i > 0$ for $i = 1, \dots, n$, and from (89), the maximizing \mathbf{S} is unique and is given by $\mathbf{S} = \mathbf{1}$. Correspondingly, from (20) we obtain the maximizing \mathbf{C}

$$\mathbf{C} = \mathbf{V}\mathbf{U}^T, \quad (90)$$

which is also unique.

Case 2: $\text{rank}[\mathbf{B}] < n$

Let $\text{rank}[\mathbf{B}] = m \leq n$. Then, $\sigma_i > 0$ for $i = 1, \dots, m$ and $\sigma_i = 0$ for $i = m + 1, \dots, n$. By (20) and (89), all maximizing solutions of Problem 6 are generated by

$$\mathbf{C} = \mathbf{V} \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{22} \end{bmatrix} \mathbf{U}^T, \quad (91)$$

where

$$\mathbf{S}_{22} \in \left\{ \mathbf{S}_{22} \in \mathbb{R}^{(n-m) \times (n-m)} : \|\mathbf{S}_{22}\| \leq 1 \right\}.$$

Clearly, the maximizing solution is non-unique.

We now summarize our findings in a Theorem.

Theorem 4 *Any local maximum of Problem 6 is a global maximum. Furthermore, Problem 6 has a unique global maximum if and only if $\text{rank}[\mathbf{B}] = n$. Let \mathbf{V} and \mathbf{U} be part of a singular value decomposition of \mathbf{B} such that*

$$\mathbf{B} = \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T,$$

where $\mathbf{V}^T\mathbf{V} = \mathbf{1}$, $\mathbf{U}^T\mathbf{U} = \mathbf{1}$ and $\mathbf{\Sigma} = \text{diag}_{i=1, \dots, n} \{\sigma_i\}$ with $\sigma_1 \geq \dots \geq \sigma_n \geq 0$. If $\text{rank}[\mathbf{B}] = n$, the unique global maximum is given by

$$\mathbf{C} = \mathbf{V}\mathbf{U}^T. \quad (92)$$

If $\text{rank}[\mathbf{B}] = m < n$, then the set of all solutions to Problem 6 is given by

$$\mathbf{C} = \mathbf{V} \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{22} \end{bmatrix} \mathbf{U}^T, \quad (93)$$

where

$$\mathbf{S}_{22} \in \left\{ \mathbf{S}_{22} \in \mathbb{R}^{(n-m) \times (n-m)} : \|\mathbf{S}_{22}\| \leq 1 \right\}.$$

Comparing Theorems 3 and 4, we see that when \mathbf{B} has full rank, Problems 5 and 6 have the same solution. Consequently (as in the comments following Theorem 3), when $\det \mathbf{B} > 0$, the solution of Problem 6 yields the solution of Problem 1. When $\det \mathbf{B} < 0$, the solution of Problem 6 yields the solution of Problem 4. In particular, the solution of Problem 2 may be found by solving Problem 6.

We have now established an equivalence between solutions of Problems 1, 2, 4, 5 with the solution of Problem 6 when \mathbf{B} has full rank. However, the solutions in all Problems are characterized by the singular value decomposition of \mathbf{B} . We seek a computational procedure for finding the solution without requiring the singular value decomposition when \mathbf{B} has full rank. To this end, we note that the norm constraint in Problem 6 is equivalent to

$$\mathbf{C}^T \mathbf{C} \leq \mathbf{1}. \quad (94)$$

Recall the Schur complement [19, ch. 2]: for $\Phi_{11} = \Phi_{11}^T \in \mathbb{R}^{p \times p}$, $\Phi_{12} \in \mathbb{R}^{p \times q}$, $\Phi_{21} \in \mathbb{R}^{q \times p}$, $\Phi_{22} = \Phi_{22}^T \in \mathbb{R}^{q \times q}$ where $\Phi_{22} > \mathbf{0}$, then

$$\Phi_{11} - \Phi_{12} \Phi_{22}^{-1} \Phi_{21} \geq \mathbf{0} \Leftrightarrow \begin{bmatrix} \Phi_{11} & \Phi_{12} \\ \Phi_{21} & \Phi_{22} \end{bmatrix} \geq \mathbf{0}.$$

Therefore, setting $\Phi_{11} = \mathbf{1}$, $\Phi_{12} = \mathbf{C}^T$, $\Phi_{21} = \mathbf{C}$, $\Phi_{22} = \mathbf{1}$ and using the Schur complement, (94) is in turn equivalent to the LMI

$$\begin{bmatrix} \mathbf{1} & \mathbf{C}^T \\ \mathbf{C} & \mathbf{1} \end{bmatrix} \geq \mathbf{0}.$$

Therefore, Problem 6 may be recast as the LMI problem

Problem 7 Let $\mathbf{B} \in \mathbb{R}^{n \times n}$. Find $\mathbf{C} \in \mathbb{R}^{n \times n}$ to

$$\text{minimize } \hat{J} = -\text{tr}[\mathbf{C}\mathbf{B}^T], \quad (95)$$

subject to

$$\begin{bmatrix} \mathbf{1} & \mathbf{C}^T \\ \mathbf{C} & \mathbf{1} \end{bmatrix} \geq \mathbf{0}. \quad (96)$$

Note that to be compatible with available LMI solvers, Problem 7 is written as a minimization problem. We now summarize our findings in the following Theorem.

Theorem 5 *Problem 7 with $\text{rank}[\mathbf{B}] = n$ has a unique global minimum, with no other local minima. When $\det[\mathbf{B}] > 0$, the solution of Problem 7 is equal to the unique solution of Problem 1. When $\det[\mathbf{B}] < 0$, the solution of Problem 7 is equal to the unique solution of Problem 4.*

We note that Problem 7 may be easily solved using existing LMI solvers. Finally, as we have seen in Theorem 4, Problem 7 has infinite solutions when $\text{rank}[\mathbf{B}] < n$, and in this case, the solution set includes matrices that are not orthonormal. Since the goal is to solve Problems 1, 2, 4, or 5, Problem 7 is useful only when \mathbf{B} has full rank (which is always the case in Problem 2).

Let us now further examine the sign of $\det[\mathbf{B}]$. Suppose that the measured vectors $\mathbf{s}_{b,k} \in \mathbb{R}^n$ are generated according to

$$\mathbf{s}_{b,k} = \bar{\mathbf{s}}_{b,k} + \mathbf{v}_k, \quad (97)$$

where

$$\bar{\mathbf{s}}_{b,k} = \mathbf{C}\mathbf{s}_{I,k}, \quad (98)$$

and $\mathbf{v}_k \in \mathbb{R}^n$ is a measurement error. Then, (5) becomes

$$\mathbf{B}^T = \bar{\mathbf{B}}^T + \Delta\mathbf{B}^T, \quad (99)$$

where

$$\bar{\mathbf{B}}^T = \sum_{k=1}^N w_k \mathbf{s}_{I,k} \bar{\mathbf{s}}_{b,k}^T, \quad \Delta\mathbf{B}^T = \sum_{k=1}^N w_k \mathbf{s}_{I,k} \mathbf{v}_k^T. \quad (100)$$

Note that $\bar{\mathbf{B}}$ is what \mathbf{B} would be if computed with error-free versions of the same vectors. We can now rewrite $\bar{\mathbf{B}}^T$ and $\Delta\mathbf{B}^T$ in (100) as

$$\bar{\mathbf{B}}^T = \mathbf{S}_I \bar{\mathbf{W}} \mathbf{S}_b^T, \quad \Delta\mathbf{B}^T = \mathbf{S}_I \bar{\mathbf{W}} \bar{\mathbf{V}}^T \quad (101)$$

where

$$\mathbf{S}_I = [\mathbf{s}_{I,1} \cdots \mathbf{s}_{I,N}], \quad \bar{\mathbf{W}} = \text{diag}\{w_1, \dots, w_N\}, \quad \mathbf{S}_b = [\bar{\mathbf{s}}_{b,1} \cdots \bar{\mathbf{s}}_{b,N}],$$

and

$$\bar{\mathbf{V}} = [\mathbf{v}_1 \cdots \mathbf{v}_N].$$

From (98), we obtain

$$\mathbf{S}_b = \mathbf{C}\mathbf{S}_I.$$

Therefore, we obtain

$$\bar{\mathbf{B}}^T = \mathbf{S}_I \bar{\mathbf{W}} \mathbf{S}_I^T \mathbf{C}^T. \quad (102)$$

Consequently, we have

$$\det \bar{\mathbf{B}} = \det (\mathbf{S}_I \bar{\mathbf{W}} \mathbf{S}_I^T) \det \mathbf{C}.$$

When $\bar{\mathbf{B}}$ has full rank, the matrix $\mathbf{S}_I \bar{\mathbf{W}} \mathbf{S}_I^T$ is positive definite, and we have

$$\text{sign}[\det \bar{\mathbf{B}}] = \text{sign}[\det \mathbf{C}]. \quad (103)$$

Finally, by continuity of the determinant and (99) together with (101), it is concluded that for a given \mathbf{S}_I and weight $\bar{\mathbf{W}}$, there exists $\delta > 0$ such that

$$\text{sign}[\det \mathbf{B}] = \text{sign}[\det \mathbf{C}], \quad \forall \bar{\mathbf{V}} \in \mathbb{R}^{n \times N} \text{ such that } \|\bar{\mathbf{V}}\| < \delta. \quad (104)$$

That is, if the collection of vectors $\mathbf{s}_{i,k}$ is geometrically rich enough, and the measurement errors \mathbf{v}_k are small enough, $\det \mathbf{B}$ will have the same sign as $\det \mathbf{C}$.

Numerical Example

In this section, Problem 1 will be solved for $n = 3, \dots, 50$ using Theorem 5, which is to say matrices $\mathbf{C} \in SO(n)$ for $n = 3, \dots, 50$ will be solved for weights w_k , vectors $\mathbf{s}_{I,k}$, and vectors $\mathbf{s}_{b,k}$, where $k = 1 \dots N$. All computations are done on a MacBook Pro with a 2.3 GHz Intel Core i5 processor with 4 GB of RAM running MATLAB 7.12.0 (R2011a). The software YALMIP [20] and SeDuMi [21], both of which interface with MATLAB in a simple and straightforward manner, are used to solve the LMI problem associated with Theorem 5 (i.e., Problem 7). The results obtained using the LMI method of Theorem 5 are compared to the SVD solution method presented in Theorem 1. All random numbers are generated using the MATLAB command `randn`.

For each n and for $k = 1 \dots N$ where $N = n + 5$, the vectors $\mathbf{s}_{I,k} \in \mathbb{R}^n$ are created by normalizing random vectors sampled from a normal distribution of mean $\mathbf{0}_n$ and covariance of $\mathbf{1}_{n \times n}$. The “noisy” measurements $\mathbf{s}_{b,k} \in \mathbb{R}^n$ are generated via $\mathbf{s}_{b,k} = \mathbf{C}'_k \mathbf{C} \mathbf{s}_{I,k}$, $k = 1, \dots, N$, where $\mathbf{C} \in SO(n)$ is the true matrix to be solved for, and $\mathbf{C}'_k \in SO(n)$ is a perturbation representing noise. The matrix \mathbf{C} is generated using a polar decomposition. Specifically, a matrix \mathbf{A} is first generated where each element of \mathbf{A} is sampled from a normal distribution with mean zero and standard deviation of one. Using a polar decomposition, the matrix \mathbf{A} can be written

$$\mathbf{A} = \mathbf{U}\mathbf{H},$$

where \mathbf{U} is orthonormal and \mathbf{H} is positive semidefinite. After checking that $\det[\mathbf{U}] = +1$ (and if it is not, a new \mathbf{A} is generated) the matrix \mathbf{C} is set equal to \mathbf{U} , that is $\mathbf{C} = \mathbf{U}$. To generate \mathbf{C}'_k , $k = 1, \dots, N$, a matrix $\mathbf{D}_k = \mathbf{1} + \mathbf{A}_k$ is first created, where each element of \mathbf{A}_k is sampled from a normal distribution of mean zero and standard deviation $1/100$. The rank of \mathbf{D}_k is checked to make sure it is full rank. Given \mathbf{D}_k the matrix \mathbf{C}'_k is found by solving Problem 2. In order to construct \mathbf{B} in (5), the weights w_k are set equal to $w_k = 1/\sigma^2$ where $\sigma = 1/100$. The weights are set to the same value because the noise on each measurement is generated the same way. Each \mathbf{B} generated is checked for $\text{rank}[\mathbf{B}] = n$ and $\det[\mathbf{B}] > 0$.

Let \mathbf{C}_{LMI} denote the best estimate of \mathbf{C} found by solving the LMI problem associated with Theorem 5 (i.e., Problem 7), and \mathbf{C}_{SVD} denote the best estimate of \mathbf{C} found using the SVD method associated with Theorem 1 using the MATLAB command `svd`. The accuracy of the solutions will be assessed by the values of $\sigma_{SVD} \triangleq \sigma_{\max}(\mathbf{C}_{SVD} - \mathbf{C})$ and $\sigma_{LMI} \triangleq \sigma_{\max}(\mathbf{C}_{LMI} - \mathbf{C})$ where, as before, σ_{\max} denotes the maximum singular value. Note that $\sigma_{\max}(\mathbf{C}_{SVD} - \mathbf{C})$ and $\sigma_{\max}(\mathbf{C}_{LMI} - \mathbf{C})$ bound the Euclidean norm of the error between $\bar{\mathbf{s}}_{b,k}$ and $\mathbf{s}_{b,k}$ where $\bar{\mathbf{s}}_{b,k} = \mathbf{C}_x \mathbf{s}_{I,k}$ and $\mathbf{s}_{b,k} = \mathbf{C}_x \mathbf{s}_{I,k}$ and \mathbf{C}_x represents either \mathbf{C}_{LMI} or \mathbf{C}_{SVD} , namely:

$$\begin{aligned} \|\bar{\mathbf{s}}_{b,k} - \mathbf{s}_{b,k}\|_2 &= \|\mathbf{C} \mathbf{s}_{I,k} - \mathbf{C}_{bI} \mathbf{s}_{I,k}\|_2 \\ &\leq \sigma_{\max}(\mathbf{C} - \mathbf{C}_{bI}) \|\mathbf{s}_{I,k}\|_2 \\ &= \sigma_{\max}(\mathbf{C} - \mathbf{C}_{bI}). \end{aligned}$$

For further comparison, the relative error minus one,

$$e_{rel-1} = \frac{\sigma_{LMI}}{\sigma_{SVD}} - 1,$$

and the relative time,

$$t_{rel} = \frac{t_{LMI}}{t_{SVD}},$$

are computed, where t_{LMI} and t_{SVD} are the execution times of the LMI and SVD methods, respectively, computed using MATLAB's `tic` and `toc` commands. The relative error minus one captures how close the solutions found using the LMI method and the SVD method are. If both methods produce a \mathbf{C} that is almost the same, then the ratio $\sigma_{LMI}/\sigma_{SVD}$ will be almost equal to one; one is subtracted from this ratio in order to better see just how close the two solutions are.

As the durations of computations in MATLAB can vary, for each k , and using $\mathbf{s}_{I,k}$, \mathbf{C} , \mathbf{C}'_k , and $\mathbf{s}_{b,k}$ for that k , e_{rel-1} and t_{rel} results are averaged over five runs.

Numerical results are shown in Table 1. Notice that the relative error minus one is essentially zero for all runs indicating that the LMI method and the SVD method are computing the same \mathbf{C} . Additionally, it is clear that as n increases the LMI method takes substantially longer to find \mathbf{C}_{LMI} as compared to the SVD method computing \mathbf{C}_{SVD} . This is not surprising given that the LMI tools, YALMIP and SeDuMi, are very general LMI tools that can be used to solve a variety of problems. In particular, SeDuMi is an optimization code capable of solving optimization problems over symmetric cones (a special case of semidefinite programming) with linear, quadratic, and semidefinite constraints, is able to solve problems with complex valued data, and at all times exploits sparsity [21]. It is also able to interface with other software packages and provide a certificate of infeasibility should a problem be infeasible. These additional features, although useful in general, are not needed to solve Problem 7. For example, Problem 7 does not have any quadratic or semidefinite constraints, the data used to solve Problem 7 is not complex valued, and $rank[\mathbf{B}] = n$ and $\det[\mathbf{B}] > 0$ so a solution is guaranteed to exist. It is expected that a solver customized to solve Problem 7 would be faster, but designing, implementing, and testing such as solver (e.g., a customized interior point method [18]) is beyond the scope of this work. Such a custom code would not incorporate some of the general aspect of YALMIP and SeDuMi, such as the ability to handle complex valued data and feasibility checks. However, a custom code tailored to solve only Problem 7 is not as attractive as a slightly more general code (so that it is faster than, say, YALMIP and SeDuMi) that is able to solve related navigation, guidance, and control problems. For instance, the constraint (96) is exploited in [23] for control purposes. In particular, [23] considers spacecraft attitude control where the spacecraft is required to avoid certain attitudes described by "keep out" and "keep in" zones. Other applications of the constraint (96) will be investigated in the future.

n	σ_{SVD}	σ_{LMI}	t_{SVD}	t_{LMI}	e_{rel-1}	t_{rel}
3	1.7940E-03	1.7940E-03	3.5400E-04	6.6026E-02	8.6819E-14	3.3643E+03
4	7.9687E-03	7.9687E-03	2.7048E-04	7.1014E-02	2.1516E-13	3.1173E+03
5	1.3772E-02	1.3772E-02	2.5678E-04	6.4836E-02	6.3949E-14	1.8471E+03
6	1.4967E-02	1.4967E-02	3.1977E-04	7.1435E-02	2.6645E-15	2.6403E+03
7	1.8392E-02	1.8392E-02	3.1400E-04	9.2804E-02	2.1538E-13	2.4254E+03
8	1.8332E-02	1.8332E-02	4.0175E-04	8.3916E-02	2.7600E-13	1.2025E+03
9	2.9556E-02	2.9556E-02	3.3930E-04	8.0396E-02	5.6288E-13	2.1055E+03
10	3.2965E-02	3.2965E-02	3.7606E-04	8.3450E-02	1.6365E-13	1.9593E+03
11	3.9148E-02	3.9148E-02	4.3219E-04	1.1147E-01	1.9473E-13	1.2160E+03
12	2.9790E-02	2.9790E-02	3.7970E-04	9.2983E-02	1.2279E-13	1.7104E+03
13	3.4809E-02	3.4809E-02	3.8328E-04	9.5668E-02	3.6859E-13	1.0707E+03
14	3.7397E-02	3.7397E-02	4.1163E-04	1.0200E-01	1.0650E-11	1.5014E+03
15	4.1795E-02	4.1795E-02	4.1306E-04	1.1527E-01	3.2907E-13	9.6795E+02
16	3.8511E-02	3.8511E-02	3.9282E-04	1.5481E-01	2.3048E-13	6.9498E+02
17	4.8311E-02	4.8311E-02	5.0566E-04	1.5965E-01	1.5108E-12	1.2055E+03
18	4.4466E-02	4.4466E-02	5.2556E-04	1.8546E-01	2.5779E-12	6.1290E+02
19	4.8543E-02	4.8543E-02	1.1332E-03	2.2211E-01	-7.5051E-14	5.1471E+02
20	4.7788E-02	4.7788E-02	5.5557E-04	2.2468E-01	7.4851E-13	9.2271E+02
21	6.0810E-02	6.0810E-02	5.6706E-04	2.9080E-01	6.6902E-13	7.1087E+02
22	7.1111E-02	7.1111E-02	5.5266E-04	3.4632E-01	2.6372E-12	8.3876E+02
23	6.0281E-02	6.0281E-02	5.9784E-04	4.2342E-01	8.1779E-13	6.9720E+02
24	8.1926E-02	8.1926E-02	6.1987E-04	4.9376E-01	3.2996E-13	5.8932E+02
25	6.1716E-02	6.1716E-02	6.5196E-04	5.7117E-01	1.0800E-12	8.2879E+02
26	7.7274E-02	7.7274E-02	8.1499E-04	7.2324E-01	7.3099E-12	6.8803E+02
27	8.0276E-02	8.0276E-02	7.3282E-04	9.5768E-01	1.3278E-13	8.0028E+02
28	1.1513E-01	1.1513E-01	8.1080E-04	8.9162E-01	1.7937E-12	9.1420E+02
29	8.0597E-02	8.0597E-02	7.8456E-04	1.0552E+00	4.4409E-13	1.0333E+03
30	9.6173E-02	9.6173E-02	8.6279E-04	1.2251E+00	4.9871E-13	1.0536E+03
31	1.1045E-01	1.1045E-01	8.9968E-04	1.8339E+00	6.7057E-14	1.2690E+03
32	9.3807E-02	9.3807E-02	7.8299E-04	1.8676E+00	1.5679E-12	1.4317E+03
33	8.5751E-02	8.5751E-02	8.9804E-04	2.3361E+00	4.7740E-14	1.5616E+03
34	9.2137E-02	9.2137E-02	9.6468E-04	2.4154E+00	7.7027E-13	1.7237E+03
35	1.1193E-01	1.1193E-01	9.6003E-04	3.2841E+00	1.2501E-13	2.0442E+03
36	9.3975E-02	9.3975E-02	1.0014E-03	3.7993E+00	5.3291E-15	2.3689E+03
37	1.0672E-01	1.0672E-01	1.0853E-03	4.8494E+00	3.7836E-13	2.7134E+03
38	1.0102E-01	1.0102E-01	1.0835E-03	4.8340E+00	1.6354E-12	3.0999E+03
39	1.1762E-01	1.1762E-01	1.2318E-03	5.2520E+00	2.9032E-12	3.2987E+03
40	1.2254E-01	1.2254E-01	1.2052E-03	7.1207E+00	2.7311E-14	3.7876E+03
41	1.1320E-01	1.1320E-01	1.1506E-03	7.9904E+00	9.7033E-14	4.3154E+03
42	1.2114E-01	1.2114E-01	1.3220E-03	7.9323E+00	1.5228E-12	4.5852E+03
43	1.2781E-01	1.2781E-01	1.2053E-03	9.4985E+00	1.2217E-12	5.7194E+03
44	1.5351E-01	1.5351E-01	1.4425E-03	1.2511E+01	-2.7089E-14	6.2369E+03
45	1.1414E-01	1.1414E-01	1.2257E-03	1.4281E+01	8.9928E-14	7.0234E+03
46	1.3442E-01	1.3442E-01	1.4735E-03	1.7777E+01	1.9074E-13	7.9236E+03
47	1.6134E-01	1.6134E-01	1.5070E-03	2.0667E+01	-9.2149E-12	9.0356E+03
48	1.5714E-01	1.5714E-01	1.8086E-03	2.1516E+01	1.8918E-13	9.1505E+03
49	1.4772E-01	1.4772E-01	1.4764E-03	2.3565E+01	2.8910E-13	1.0444E+04
50	1.3152E-01	1.3152E-01	1.6183E-03	2.6558E+01	1.3434E-13	1.1739E+04

Table 1 Numerical results of finding C_{LMI} and C_{SVD} for $n = 1, \dots, 50$.

Conclusion

In this paper, we have presented a rigorous analysis of the famous Wahba problem on $SO(n)$. In particular, we obtain the entire set of solutions, using

both singular value decomposition matrix square-root methods. In doing so, we also obtain conditions for uniqueness of solutions, correcting some errors in the existing literature. We then show that under a mild condition, which holds if the measured vectors are geometrically rich enough and the measurement errors are small enough, Wahba's problem on $SO(n)$ can be recast as a linear matrix inequality optimization problem, by appropriately relaxing the $SO(n)$ constraint. This opens the door to a whole host of new possible solution methods for these types of Wahba problems. It also suggests an approach for how other optimization problems on $SO(n)$ might be relaxed and made more readily solvable.

Acknowledgements The authors gratefully acknowledge discussions with Prof. Tim Barfoot who pointed out how to decouple the simultaneous translation and rotation problem on $SE(n)$. The authors also express sincere gratitude to the anonymous reviewers whose comments helped to greatly improve this manuscript.

References

1. G. Wahba, "A Least-Squares Estimate of Satellite Attitude," *SIAM Review*, 7 (3), 1965, pp. 409.
2. J. Keat, "Analysis of Least-Squares Attitude Determination Routine DOAOP," Computer Sciences Corporation Report CSC/TM-77/6034, February 1977.
3. B.K.P. Horn, "Closed-form solution of absolute orientation using unit quaternions," *Journal of the Optical Society of America*, 4, 1987, pp. 629–642.
4. M.D. Shuster, S.D. Oh, "Three-Axis Attitude Determination from Vector Observations," *AIAA Journal of Guidance, Control and Dynamics*, 4 (1), 1981, pp. 70–77.
5. F.L. Markley, D. Mortari, "Quaternion Attitude Estimation Using Vector Observations," *Journal of the Astronautical Sciences*, 48 (2-3), 2000, pp. 359–280.
6. B.F. Green, "The Orthogonal Approximation of an Oblique Structure in Factor Analysis," *Psychometrika*, 17 (4), 1952, pp. 429 – 440.
7. P.H. Schöneman, "A Generalized Solution to the Orthogonal Procrustes Problem," *Psychometrika*, 31 (1), 1966, pp. 1 – 10.
8. J.L. Farrell, J.C. Stuelpnagel, R.H. Wessner, J.R. Velman, J.E. Brook, "A Least-Squares Estimate of Satellite Attitude," *SIAM Review*, 8 (3), 1966, pp. 384–386. doi: 10.1137/1008080
9. K.S. Arun, T.S. Huang, S.D. Blostein, "Least-Squares Fitting of Two 3-D Point Sets," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 9 (5), 1987, pp. 698–700.
10. F.L. Markley, "Attitude Determination Using Vector Observations and the Singular Value Decomposition," *Journal of the Astronautical Sciences*, 36 (3), July-September 1988, pp. 245-258.
11. S. Umeyama, "Least-Squares Estimation of Transformation Parameters Between Two Point Patterns," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13 (4), 1991, pp. 376–380.
12. K. Kanatani, "Analysis of 3-D Rotation Fitting," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16 (5), 1994, pp. 543–549.
13. M.L. Psiaki, "Generalized Wahba Problems for Spinning Spacecraft Attitude and Rate Determination," *Journal of the Astronautical Sciences*, 57 (1,2), 2010, pp. 73–92.
14. J.C. Hinks, M.L. Psiaki, "Solution Strategies for an Extension of Wahbas Problem to a Spinning Spacecraft," *Journal of Guidance, Control and Dynamics*, 34 (6), 2011, pp. 1734–1745.
15. M.L. Psiaki, J.C. Hinks, "Numerical Solution of a Generalized Wahba Problem for a Spinning Spacecraft," *Journal of Guidance, Control and Dynamics*, 35 (3), 2012, pp. 764–773.

16. J. Forbes and A. de Ruiter, "An LMI-Based Solution to Wahba's Problem," *AIAA Journal of Guidance, Control and Dynamics*, to appear.
17. P.C. Hughes, *Spacecraft Attitude Dynamics*, Dover Publications, New York, 2004.
18. J. Nocedal, S.J. Wright, *Numerical Optimization*, Springer-Verlag, New York, 1999.
19. S. Boyd, L. El Ghaoui, E. Feron, V. Balakrishnan *Linear Matrix Inequalities in System and Control Theory*, Volume 15 of Studies in Applied Mathematics, Society for Industrial and Applied Mathematics (SIAM), 1994.
20. J. Loftberg, "YALMIP: A Toolbox for Modeling and Optimization in MATLAB", in CACSD Conf. Taipei, Taiwan, September 4, 2004, pp. 284 – 289.
21. J. Sturm, "Using SeDuMi 1.02, a MATLAB Toolbox for Optimization Over Symmetric Cones, Optimization Methods and Software, vol. 11, no. 12, pp. 625–653, 1999, special Issue on Interior Point Methods.
22. J. Mattingley and S. Boyd, "Real-time Convex Optimization in Signal Processing," *IEEE Signal Processing Magazine*, 27 (3), 2010, pp. 50–56.
doi: 10.1109/MSP.2010.936020
23. A. Weiss, F. Leve, M. Baldwin, J. R. Forbes, I. Kolmanovsky, "Spacecraft Constrained Attitude Control using Positively Invariant Constraint Admissible Sets on $SO(3) \times \mathbb{R}^3$," to appear in Proc. of the American Control Conference, Portland, OR, June 4 - 6, 2014.
24. J.F. Franklin, *Matrix Theory*, Dover Publications, New York, 2000.

Useful Matrix Properties

We list some basic matrix properties which are used throughout the paper. The first three lemmas present some well-known properties of real symmetric matrices [24, pp. 99–102].

Lemma 2 *Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be a real symmetric matrix. Then, all eigenvalues of \mathbf{A} are real. Furthermore, all associated eigenvectors may be taken to be real also.*

Lemma 3 *Let λ_1 and λ_2 be two distinct eigenvalues of a real symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, and let $\mathbf{x}_1 \in \mathbb{R}^n$ and $\mathbf{x}_2 \in \mathbb{R}^n$ be associated eigenvectors (that is $\mathbf{A}\mathbf{x}_1 = \lambda_1\mathbf{x}_1$ and $\mathbf{A}\mathbf{x}_2 = \lambda_2\mathbf{x}_2$). Then, \mathbf{x}_1 and \mathbf{x}_2 are orthogonal.*

Lemma 4 *Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be a real symmetric matrix. Then, there exists a real orthonormal matrix $\mathbf{V} \in \mathbb{R}^{n \times n}$ (that is $\mathbf{V}^T\mathbf{V} = \mathbf{1}$) and a real diagonal matrix $\mathbf{\Lambda} = \text{diag}\{\lambda_1, \dots, \lambda_n\}$, such that*

$$\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T,$$

where $\lambda_i \in \mathbb{R}$ for $i = 1, \dots, n$ are the eigenvalues of \mathbf{A} , and the columns of \mathbf{V} are associated eigenvectors (that is, $\mathbf{A}\mathbf{v}_i = \lambda_i\mathbf{v}_i$ for $i = 1, \dots, n$, with $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_n]$).

The following lemma characterizes all orthonormal real eigenmatrices of a real symmetric matrix, for a given ordering of its eigenvalues.

Lemma 5 *Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be a real symmetric matrix with $m \leq n$ distinct eigenvalues denoted by $\lambda_1, \dots, \lambda_m \in \mathbb{R}$ (eigenvalues may be repeated). Let n_i be the multiplicity of λ_i , for $i = 1, \dots, m$, such that $\sum_{i=1}^m n_i = n$. In addition, let $\mathbf{V}_0 \in \mathbb{R}^{n \times n}$ be a real orthonormal matrix such that*

$$\mathbf{A} = \mathbf{V}_0\mathbf{\Lambda}\mathbf{V}_0^T,$$

where

$$\mathbf{\Lambda} = \text{diag}_{i=1, \dots, m} \{ \lambda_i \mathbf{1}_{n_i \times n_i} \}.$$

Then, $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$ with orthonormal $\mathbf{V} \in \mathbb{R}^{n \times n}$ if and only if $\mathbf{V} = \mathbf{V}_0\mathbf{\Delta}$, where

$$\mathbf{\Delta} = \text{diag}_{i=1, \dots, m} \{ \mathbf{\Delta}_i \},$$

and $\mathbf{\Delta}_i \in \mathbb{R}^{n_i \times n_i}$ are real orthonormal matrices.

Proof First it will be shown that $\mathbf{V} = \mathbf{V}_0\mathbf{\Delta}$ where $\mathbf{\Delta}$ is as in the statement of the Lemma implies that $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T = \mathbf{V}_0\mathbf{\Lambda}\mathbf{V}_0^T$. Let \mathbf{V}_0 be given (by Lemma 4 and a suitable re-ordering of the columns of \mathbf{V}_0 and diagonal entries of $\mathbf{\Lambda}$, it exists), and let $\mathbf{V} = \mathbf{V}_0\mathbf{\Delta}$ as in the statement of the Lemma. It is readily apparent that \mathbf{V} is orthonormal. It follows that

$$\mathbf{V}\mathbf{\Lambda}\mathbf{V}^T = \mathbf{V}_0\mathbf{\Delta}\mathbf{\Lambda}\mathbf{\Delta}^T\mathbf{V}_0^T,$$

Next,

$$\begin{aligned}\mathbf{\Delta}\mathbf{\Lambda}\mathbf{\Delta}^T &= \text{diag}_{i=1,\dots,m} \{\mathbf{\Delta}_i\} \text{diag}_{i=1,\dots,m} \{\lambda_i \mathbf{1}_{n_i \times n_i}\} \text{diag}_{i=1,\dots,m} \{\mathbf{\Delta}_i^T\}, \\ &= \text{diag}_{i=1,\dots,m} \{\lambda_i \mathbf{\Delta}_i \mathbf{\Delta}_i^T\}, \\ &= \text{diag}_{i=1,\dots,m} \{\lambda_i \mathbf{1}_{n_i \times n_i}\}, \\ &= \mathbf{\Lambda}.\end{aligned}$$

Hence, it follows that

$$\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T = \mathbf{V}_0\mathbf{\Lambda}\mathbf{V}_0^T.$$

Now it will be shown that \mathbf{V} and \mathbf{V}_0 satisfying $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T = \mathbf{V}_0\mathbf{\Lambda}\mathbf{V}_0^T$ implies that $\mathbf{\Delta}$ is as in the statement of the Lemma. Suppose that $\mathbf{V} \in \mathbb{R}^{n \times n}$ and $\mathbf{V}_0 \in \mathbb{R}^{n \times n}$ are orthonormal matrices satisfying

$$\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T = \mathbf{V}_0\mathbf{\Lambda}\mathbf{V}_0^T.$$

Let us now partition both \mathbf{V}_0 and \mathbf{V} as

$$\mathbf{V}_0 = \text{col}_{i=1,\dots,m} \{\mathbf{V}_{0,i}\}, \quad \mathbf{V} = \text{col}_{i=1,\dots,m} \{\mathbf{V}_i\},$$

where

$$\mathbf{V}_{0,i} = \text{col}_{j=1,\dots,n_i} \{\mathbf{v}_{0,i}^j\}, \quad \mathbf{V}_i = \text{col}_{j=1,\dots,n_i} \{\mathbf{v}_i^j\}, \quad \text{for } i = 1, \dots, m.$$

That is, each $\mathbf{V}_{0,i}$ and \mathbf{V}_i contains eigenvectors of \mathbf{A} corresponding to eigenvalue λ_i , for $i = 1, \dots, m$. Now, let $\mathbf{x} \in \mathbb{R}^n$ be any real eigenvector of \mathbf{A} corresponding to the eigenvalue λ_i , for some $i \in \{1, \dots, m\}$. Since \mathbf{V}_0 is orthonormal, its columns span \mathbb{R}^n , and \mathbf{x} can be written as a linear combination of them. However, by Lemma 3, \mathbf{x} is orthogonal to all columns of $\mathbf{V}_{0,k}$, for $k = 1, \dots, m$ with $k \neq i$. Therefore, it follows that \mathbf{x} can be written as a linear combination of the columns of $\mathbf{V}_{0,i}$. Since the columns of \mathbf{V}_i are all eigenvectors of \mathbf{A} corresponding to the eigenvalue λ_i , it follows that \mathbf{V}_i may be written as

$$\mathbf{V}_i = \mathbf{V}_{0,i}\mathbf{\Delta}_i,$$

for some $\mathbf{\Delta}_i \in \mathbb{R}^{n_i \times n_i}$. Since \mathbf{V} is orthonormal, it follows that $\mathbf{V}_i^T \mathbf{V}_i = \mathbf{1}_{n_i \times n_i}$, from which we obtain

$$\begin{aligned}\mathbf{1}_{n_i \times n_i} &= \mathbf{V}_i^T \mathbf{V}_i, \\ &= \mathbf{\Delta}_i^T \mathbf{V}_{0,i}^T \mathbf{V}_{0,i} \mathbf{\Delta}_i, \\ &= \mathbf{\Delta}_i^T \mathbf{\Delta}_i.\end{aligned}$$

Hence, $\mathbf{\Delta}_i$ is orthonormal also.

The next lemma builds on the previous one, to characterize all singular value decompositions of a real square matrix.

Lemma 6 *Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be a real square matrix with singular values $\sigma_1 > \dots > \sigma_m \geq 0$ with $m \leq n$ (singular values may be repeated). Let n_i be the multiplicity of σ_i , for $i = 1, \dots, m$, such that $\sum_{i=1}^m n_i = n$. In addition, let $\mathbf{V}_0 \in \mathbb{R}^{n \times n}$ and $\mathbf{U}_0 \in \mathbb{R}^{n \times n}$ be real orthonormal matrices such that \mathbf{A} has a singular value decomposition*

$$\mathbf{A} = \mathbf{V}_0 \mathbf{\Sigma} \mathbf{U}_0^T,$$

where

$$\boldsymbol{\Sigma} = \text{diag} \left\{ \sigma_i \mathbf{1}_{n_i \times n_i} \right\}_{i=1, \dots, m}.$$

Then, $\mathbf{A} = \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T$ with orthonormal $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{n \times n}$ if and only if $\mathbf{V} = \mathbf{V}_0 \boldsymbol{\Delta}$ and $\mathbf{U} = \mathbf{U}_0 \bar{\boldsymbol{\Delta}}$, where

$$\boldsymbol{\Delta} = \text{diag} \left\{ \boldsymbol{\Delta}_i \right\}_{i=1, \dots, m}, \quad \bar{\boldsymbol{\Delta}} = \text{diag} \left\{ \bar{\boldsymbol{\Delta}}_i \right\}_{i=1, \dots, m},$$

and $\boldsymbol{\Delta}_i, \bar{\boldsymbol{\Delta}}_i \in \mathbb{R}^{n_i \times n_i}$ are real orthonormal matrices. Furthermore,

$$\boldsymbol{\Delta}_i = \bar{\boldsymbol{\Delta}}_i \text{ for } i = 1, \dots, m-1.$$

If $\sigma_m > 0$ (\mathbf{A} has full rank), then

$$\boldsymbol{\Delta}_m = \bar{\boldsymbol{\Delta}}_m,$$

in which case $\boldsymbol{\Delta} = \bar{\boldsymbol{\Delta}}$.

If $\sigma_m = 0$ (\mathbf{A} is rank deficient), then $\boldsymbol{\Delta}_m$ and $\bar{\boldsymbol{\Delta}}_m$ may be chosen independently, and $\boldsymbol{\Delta}$ and $\bar{\boldsymbol{\Delta}}$ are not necessarily equal.

Proof First it will be shown that $\mathbf{V} = \mathbf{V}_0 \boldsymbol{\Delta}$ and $\mathbf{U} = \mathbf{U}_0 \bar{\boldsymbol{\Delta}}$ where $\boldsymbol{\Delta}$ and $\bar{\boldsymbol{\Delta}}$ are as in the statement of the Lemma implies that $\mathbf{A} = \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T = \mathbf{V}_0 \boldsymbol{\Sigma} \mathbf{U}_0^T$. To begin, let $\boldsymbol{\Delta}$ and $\bar{\boldsymbol{\Delta}}$ be given as in the Lemma. It is readily apparent that \mathbf{V} and \mathbf{U} are orthonormal. Next,

$$\mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T = \mathbf{V}_0 \boldsymbol{\Delta} \boldsymbol{\Sigma} \bar{\boldsymbol{\Delta}}^T \mathbf{U}_0^T,$$

$$\begin{aligned} \boldsymbol{\Delta} \boldsymbol{\Sigma} \bar{\boldsymbol{\Delta}}^T &= \text{diag} \left\{ \text{diag} \left\{ \boldsymbol{\Delta}_i \right\}_{i=1, \dots, m-1}, \boldsymbol{\Delta}_m \right\} \text{diag} \left\{ \text{diag} \left\{ \sigma_i \mathbf{1}_{n_i \times n_i} \right\}_{i=1, \dots, m-1}, \sigma_m \mathbf{1}_{n_m \times n_m} \right\} \\ &\quad \times \text{diag} \left\{ \text{diag} \left\{ \boldsymbol{\Delta}_i^T, \bar{\boldsymbol{\Delta}}_m^T \right\}_{i=1, \dots, m-1} \right\}, \\ &= \text{diag} \left\{ \text{diag} \left\{ \sigma_i \boldsymbol{\Delta}_i \boldsymbol{\Delta}_i^T, \sigma_m \boldsymbol{\Delta}_m \bar{\boldsymbol{\Delta}}_m^T \right\}_{i=1, \dots, m-1} \right\}, \\ &= \text{diag} \left\{ \text{diag} \left\{ \sigma_i \mathbf{1}_{n_i \times n_i}, \sigma_m \boldsymbol{\Delta}_m \bar{\boldsymbol{\Delta}}_m^T \right\}_{i=1, \dots, m-1} \right\}. \end{aligned}$$

Now, when $\sigma_m > 0$, we require $\boldsymbol{\Delta}_m = \bar{\boldsymbol{\Delta}}_m$, and

$$\sigma_m \boldsymbol{\Delta}_m \bar{\boldsymbol{\Delta}}_m^T = \sigma_m \mathbf{1}_{n_m \times n_m}.$$

When $\sigma_m = 0$, we have

$$\sigma_m \boldsymbol{\Delta}_m \bar{\boldsymbol{\Delta}}_m^T = \mathbf{0}_{n_m \times n_m} = \sigma_m \mathbf{1}_{n_m \times n_m}.$$

In both cases, we obtain

$$\boldsymbol{\Delta} \boldsymbol{\Sigma} \bar{\boldsymbol{\Delta}}^T = \boldsymbol{\Sigma}.$$

Hence, it follows that

$$\mathbf{A} = \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T = \mathbf{V}_0 \boldsymbol{\Sigma} \mathbf{U}_0^T.$$

Now it will be shown that $\mathbf{V}, \mathbf{U}, \mathbf{V}_0$, and \mathbf{U}_0 satisfying $\mathbf{A} = \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T = \mathbf{V}_0 \boldsymbol{\Sigma} \mathbf{U}_0^T$ implies that $\boldsymbol{\Delta}$ and $\bar{\boldsymbol{\Delta}}$ are as in the statement of the Lemma.

Let $\mathbf{V}, \mathbf{U}, \mathbf{V}_0$, and \mathbf{U}_0 be any real orthonormal matrices satisfying

$$\mathbf{A} = \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T = \mathbf{V}_0 \boldsymbol{\Sigma} \mathbf{U}_0^T. \quad (105)$$

We then readily find that

$$\mathbf{A} \mathbf{A}^T = \mathbf{V} \boldsymbol{\Sigma}^2 \mathbf{V}^T, \quad \mathbf{A}^T \mathbf{A} = \mathbf{U} \boldsymbol{\Sigma}^2 \mathbf{U}^T. \quad (106)$$

Note that this is in particular true for \mathbf{V}_0 and \mathbf{U}_0 . Now, if we let $\mathbf{A}' = \mathbf{A}\mathbf{A}^T = \mathbf{V}\boldsymbol{\Sigma}^2\mathbf{V}^T$ and $\mathbf{A}'' = \mathbf{A}^T\mathbf{A} = \mathbf{U}\boldsymbol{\Sigma}^2\mathbf{U}^T$, notice the similarity between the structure of \mathbf{A}' and \mathbf{A}'' to the structure used in Lemma 5. As such, by Lemma 5, the sets of all real orthonormal \mathbf{V} and \mathbf{U} satisfying (106) are given by

$$\mathbf{V} = \mathbf{V}_0\boldsymbol{\Delta}, \quad \mathbf{U} = \mathbf{U}_0\bar{\boldsymbol{\Delta}},$$

where

$$\boldsymbol{\Delta} = \text{diag}_{i=1,\dots,m} \{\boldsymbol{\Delta}_i\}, \quad \bar{\boldsymbol{\Delta}} = \text{diag}_{i=1,\dots,m} \{\bar{\boldsymbol{\Delta}}_i\},$$

and $\boldsymbol{\Delta}_i, \bar{\boldsymbol{\Delta}}_i \in \mathbb{R}^{n_i \times n_i}$ are real orthonormal matrices. Therefore, (105) can be written as

$$\mathbf{A} = \mathbf{V}_0\boldsymbol{\Delta}\boldsymbol{\Sigma}\bar{\boldsymbol{\Delta}}^T\mathbf{U}_0^T. \quad (107)$$

Rearranging (107), we find that

$$\boldsymbol{\Delta}\boldsymbol{\Sigma}\bar{\boldsymbol{\Delta}}^T = \mathbf{V}_0^T\mathbf{A}\mathbf{U}_0. \quad (108)$$

However, from (105) where $\mathbf{A} = \mathbf{V}_0\boldsymbol{\Sigma}\mathbf{U}_0^T$, we also have

$$\mathbf{V}_0^T\mathbf{A}\mathbf{U}_0 = \boldsymbol{\Sigma}. \quad (109)$$

Hence, comparing (108) and (109), we have

$$\boldsymbol{\Delta}\boldsymbol{\Sigma}\bar{\boldsymbol{\Delta}}^T = \boldsymbol{\Sigma}.$$

Expanding, this becomes

$$\text{diag}_{i=1,\dots,m} \{\boldsymbol{\Delta}_i\} \text{diag}_{i=1,\dots,m} \{\sigma_i \mathbf{1}_{n_i \times n_i}\} \text{diag}_{i=1,\dots,m} \{\bar{\boldsymbol{\Delta}}_i^T\} = \text{diag}_{i=1,\dots,m} \{\sigma_i \mathbf{1}_{n_i \times n_i}\},$$

from which we obtain

$$\sigma_i \boldsymbol{\Delta}_i \bar{\boldsymbol{\Delta}}_i^T = \sigma_i \mathbf{1}_{n_i \times n_i}, \quad i = 1, \dots, m. \quad (110)$$

Since $\sigma_i > 0$ for $i = 1, \dots, m-1$, this immediately leads to $\boldsymbol{\Delta}_i \bar{\boldsymbol{\Delta}}_i^T = \mathbf{1}_{n_i \times n_i}$. Post-multiplying both sides by $\boldsymbol{\Delta}_i$ gives

$$\boldsymbol{\Delta}_i = \bar{\boldsymbol{\Delta}}_i, \quad i = 1, \dots, m-1.$$

Now, for $i = m$, there are two cases. In the first case, $\sigma_m > 0$ (\mathbf{A} has full rank). In this case, just as for $i = 1, \dots, m-1$, we obtain

$$\boldsymbol{\Delta}_m = \bar{\boldsymbol{\Delta}}_m.$$

In the second case, $\sigma_m = 0$ (\mathbf{A} is rank deficient), and (110) is automatically satisfied for any orthonormal $\boldsymbol{\Delta}_m$ and $\bar{\boldsymbol{\Delta}}_m$.

The following proposition and its corollary are obtained from [16].

Proposition 1 Consider any matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$, with $\|\mathbf{A}\| = \ell$, for some $\ell \geq 0$. Denote the ij^{th} term of \mathbf{A} as a_{ij} . Then,

$$\sqrt{\sum_{i=1}^n a_{ij}^2} \leq \ell, \quad j = 1, \dots, m. \quad (111)$$

and

$$\sqrt{\sum_{j=1}^m a_{ij}^2} \leq \ell, \quad i = 1, \dots, n, \quad (112)$$

Corollary 1 Consider any matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$, with $\|\mathbf{A}\| = \ell$, for some $\ell \geq 0$. Denote the ij^{th} term of \mathbf{A} as a_{ij} . Then,

$$|a_{ij}| \leq \ell.$$

Proposition 1 implies that no row or column can have a Euclidean norm greater than the induced matrix 2-norm, while Corollary 1 means that no individual matrix element can be greater than the induced matrix 2-norm.