

Some studies on finding the nearest volume-preserving matrix

Mu Jianfei Huang Jianguo

(Department of Mathematics, Shanghai Jiaotong University, Shanghai 200240, China)

Abstract: Finding the nearest volume-preserving matrix for a given matrix is studied. A matrix equation is first obtained, which is a necessary condition for the solution to the problem. Then the equation is solved by the singular value decomposition method. Some additional results are also provided to further characterize the solution. Using these results, a numerical algorithm is introduced and a numerical test is given to illustrate the effectiveness of the algorithm.

Key words: volume-preserving matrix; matrix nearness problem; singular value decomposition

We consider the following nearest volume-preserving matrix problem: For any given matrix $A \in \mathbf{R}^{n \times n}$, find a matrix $\hat{X} \in S(n)$ such that

$$\hat{X} = \arg \min_{X \in S(n)} \|A - X\|_F^2 \quad (1)$$

where $S(n)$ denotes the set of all real volume-preserving matrices of order n ,

$$S(n) \equiv \{A \in \mathbf{R}^{n \times n} : \det(A) = 1\}$$

which is a Lie group with the Lie algebra^[1]:

$$s(n) \equiv \{A \in \mathbf{R}^{n \times n} : \text{tr}(A) = 0\}$$

Such kinds of problems have wide applications in signal processing and computer vision (see Refs. [2 – 7]). The problem considered here can be viewed as a “dual” problem of the one considered in Ref. [7], in the sense that the constraint sets of the two problems are $S(n)$ and $s(n)$ respectively, and the latter one is just the Lie algebra of the former one. On the other hand, observing that $S(n)$ consists of all volume-preserving affine transformations which have practical applications^[8], when we obtain an approximate transformation by measurements, we have to solve problem (1) to get the modified one with a physical sense. We refer to Refs. [3, 6, 9] for more details along this line.

In this paper, we first derive a matrix equation which should be satisfied by the solution to problem (1). We then solve the equation by the singular value decomposition method. We have also obtained a result characterizing the solution of (1) in detail. Finally, a numerical algorithm is provided by means of the previous results and a numerical test is given to illustrate the effectiveness of the algorithm.

1 Necessary Conditions for Solutions of (1)

In this section, we first give a necessary condition

for the solution of (1).

Lemma 1 Let $\hat{X} \in S(n)$ solve the least square problem (1). Then it must satisfy the following matrix equation:

$$X^T(X - A) = \lambda I_n \quad X \in S(n) \quad (2)$$

where $\lambda \in \mathbf{R}$ is some real constant and I_n denotes the $n \times n$ identity matrix.

Proof It is evident for any matrix $W \in s(n)$ and any $t \in \mathbf{R}$ that $\det(e^{tW}) = 1$. Hence, the smooth curve $\{\hat{X}e^{tW}\}_{t \in \mathbf{R}}$ is on the Lie group $S(n)$. Noting that \hat{X} is the solution of (1), we find

$$\left. \frac{d}{dt} \|A - \hat{X}e^{tW}\|_F^2 \right|_{t=0} = 0$$

which together with some simple computation implies

$$\langle A - \hat{X}, \hat{X}W \rangle = 0$$

i. e. ,

$$\langle \hat{X}^T(A - \hat{X}), W \rangle = 0 \quad (3)$$

Moreover, since $Z - \frac{1}{n} \text{tr}(Z) I_n \in s(n)$ for any $Z \in \mathbf{R}^{n \times n}$, it follows from (3) that

$$\langle \hat{X}^T(A - \hat{X}), Z - \frac{1}{n} \text{tr}(Z) I_n \rangle = 0$$

i. e. ,

$$\langle \hat{X}^T(A - \hat{X}) - \frac{1}{n} \text{tr}(\hat{X}^T(A - \hat{X})) I_n, Z \rangle = 0$$

which yields asserted result due to the arbitrariness of Z .

We next consider the solution to the matrix equation (2). To avoid unnecessary complexity, we assume in what follows that A is nonsingular, and so admits the singular value decomposition^[10]:

$$A = P_1 \Lambda P_2 \quad (4)$$

where $P_1, P_2 \in O(n)$; $\Lambda = \text{diag}(\mu_1, \mu_2, \dots, \mu_n)$, $\mu_1 \geq \mu_2 \geq \dots \geq \mu_n > 0$; $O(n)$ denotes the set of all $n \times n$ orthogonal matrices, and $\text{diag}(\tau_1, \tau_2, \dots, \tau_n)$ is a diagonal matrix with the i -th diagonal entry τ_i .

Lemma 2 Let the $n \times n$ matrix A be nonsingular with the singular value decomposition (4). Then

Received 2004-06-02.

Biography: Mu Jianfei (1961—), male, associate professor, jfmu@sjtu.edu.cn.

there exist some constants $\tau_1, \tau_2, \dots, \tau_n$ such that

$$X = P_1 \text{diag}(\tau_1, \tau_2, \dots, \tau_n) P_2 \quad (5)$$

solves the matrix equation (2). Moreover, each solution of Eq. (2) must be described in this form by means of a proper singular value decomposition of A .

Proof Let X be a solution of Eq. (2). It follows from Eq. (2) that $X^T A$ is symmetric,

$$X^T A = A^T X \quad (6)$$

Since the singular values $\mu_1, \mu_2, \dots, \mu_n$ of A may be equal, we rewrite the singular value decomposition (4) in block form for later use.

$$A = P_1 \Lambda P_2 \quad (7)$$

where $P_1, P_2 \in O(n)$, $\Lambda = \text{diag}(\sigma_1 I_{n_1}, \sigma_2 I_{n_2}, \dots, \sigma_k I_{n_k})$, $\sigma_1 > \sigma_2 > \dots > \sigma_k > 0$.

Let $X_1 = P_1^T X P_2^T$. Then Eq. (6) can be recast as

$$(\Lambda X_1)^T = \Lambda X_1$$

which indicates $\Lambda X_1 = V$ is symmetric, hence Eq. (2) becomes

$$V \Lambda^{-2} V - V = \lambda I_n$$

or equivalently,

$$\Lambda^{-2} V = I_n + \lambda V^{-1}$$

So we know $\Lambda^{-2} V$ is symmetric, which gives

$$\Lambda^{-2} V = V \Lambda^{-2} \quad (8)$$

Now let us partition the matrix V in the following block form compatible with the block structure of Λ :

$$V = \begin{bmatrix} V_{11} & \cdots & V_{1k} \\ \vdots & & \vdots \\ V_{k1} & \cdots & V_{kk} \end{bmatrix} \quad V_{ij} \in \mathbf{R}^{n_i \times n_j}$$

Then Eq. (8) is equivalent to

$$(\sigma_i^{-2} - \sigma_j^{-2}) V_{ij} = \mathbf{0} \quad i, j = 1, 2, \dots, k$$

which leads to

$$V_{ij} = \mathbf{0} \quad i \neq j \\ V_{ii} = V_{ii}^T \in \mathbf{R}^{n_i \times n_i} \quad 1 \leq i \leq k$$

On the other hand, since V_{ii} is symmetric, we have the spectral decomposition:

$$V_{ii} = Q_{ii} \Delta_{ii} Q_{ii}^T \quad Q_{ii} \in O(n_i)$$

where Δ_{ii} is diagonal. We define

$$Q_1 \equiv P_1 \begin{bmatrix} Q_{11} & & \\ & \ddots & \\ & & Q_{kk} \end{bmatrix} \\ Q_2 \equiv \begin{bmatrix} Q_{11} & & \\ & \ddots & \\ & & Q_{kk} \end{bmatrix}^T P_2$$

and

$$\Delta \equiv \begin{bmatrix} \sigma_1^{-1} \Delta_{11} & & \\ & \ddots & \\ & & \sigma_k^{-1} \Delta_{kk} \end{bmatrix}$$

Then, it is clear that $Q_1, Q_2 \in O(n)$,

$$A = P_1 \Lambda P_2 = Q_1 \Delta Q_2$$

and

$$X = P_1 X_1 P_2 = P_1 \Lambda^{-1} V P_2 = \\ P_1 \Lambda^{-1} \begin{bmatrix} Q_{11} \Delta_{11} Q_{11}^T & & \\ & \ddots & \\ & & Q_{kk} \Delta_{kk} Q_{kk}^T \end{bmatrix} P_2 = Q_1 \Delta Q_2$$

Now, if A and X are written in the above forms, we can easily see that Eq. (2) is equivalent to $\Delta^2 - \Delta A = \lambda I_n$. Hence, $X = P_1 \Delta P_2$ also solves Eq. (2) if Δ meets the previous equation. Note also that $A = Q_1 \Lambda Q_2$ is a singular value decomposition of A . Therefore, each solution of Eq. (2) can be described in the form Eq. (5) by means of a proper singular value decomposition of A .

2 Further Discussions

It follows from lemma 1 and lemma 2 that each solution \hat{X} of the constrained optimization problem (1) can always be expressed in the form (5). Therefore, in order to solve problem (1), it suffices to evaluate all the values of $\|A - X\|_F^2$ with X given by Eq. (5). In this case, we have by the basic property of Frobenius norm that

$$g(\tau) \equiv \|A - X\|_F^2 = \|A - \text{diag}(\tau_1, \tau_2, \dots, \tau_n)\|_F^2 = \\ \sum_{i=1}^n (\mu_i - \tau_i)^2 \quad (9)$$

where, as shown in lemma 2, $\{\mu_i\}_{i=1}^n$ are the singular values of A enumerated in descending order. Moreover, the constraint condition $\det(X) = 1$ amounts to

$$\prod_{i=1}^n \tau_i = \det(Q_1 Q_2) = \text{sgn}(\det(A)) \quad (10)$$

where $\text{sgn}(\cdot)$ is a sign function defined by

$$\text{sgn}(x) = \begin{cases} 1 & x \geq 0 \\ -1 & x < 0 \end{cases}$$

Thus problem (1) can be reformulated as

$$\min_{\tau} g(\tau) \quad (11)$$

$$\text{subject to } \prod_{i=1}^n \tau_i = \text{sgn}(\det(A))$$

Lemma 3 Let $\hat{\tau} = \{\hat{\tau}_1, \hat{\tau}_2, \dots, \hat{\tau}_n\}$ solve the constrained optimization problem (11). Then there exists some real constant λ such that

$$\hat{\tau}_i = \frac{\mu_i \pm \sqrt{\mu_i^2 + \lambda}}{2} \quad i = 1, 2, \dots, n \quad (12)$$

Proof By the method of Lagrange multipliers, there exists some real constant $\frac{\text{sgn}(\det(A))}{2} \lambda$ such

that $\hat{\tau}$ satisfies

$$\partial_{\tau_i} \left[\sum_{i=1}^n (\mu_i - \tau_i)^2 - \frac{\text{sgn}(\det(A))}{2} \lambda \cdot \right.$$

$$\left. \left(\prod_{i=1}^n \tau_i - \text{sgn}(\det(A)) \right) \right] = 0 \quad i = 1, 2, \dots, n$$

or equivalently,

$$\tau_i(\mu_i - \tau_i) = -\frac{\text{sgn}(\det(\mathbf{A}))}{4} \lambda \prod_{i=1}^n \tau_i = -\frac{\lambda}{4}$$

which yields the desired expression (12) immediately.

Lemma 3 provides an important property characterizing the solution of problem (11), but there are still $2n$ possible combinations, each of which may result in the desired solution. The next result is given to simplify this difficulty.

Lemma 4 Let $\hat{\tau} = \{\hat{\tau}_1, \hat{\tau}_2, \dots, \hat{\tau}_n\}$ solve the constrained optimization problem (11). Then the following statements hold:

$$\textcircled{1} \hat{\tau}_1 \geq \hat{\tau}_2 \geq \dots \geq \hat{\tau}_n;$$

$\textcircled{2}$ If there exists some $\hat{\tau}_i$ such that $\hat{\tau}_i = \frac{\mu_i - \sqrt{\mu_i^2 + \lambda}}{2}$, then for every $j \geq i$, $\mu_j^2 + \lambda > 0$, we have

$$\hat{\tau}_j = \frac{\mu_j - \sqrt{\mu_j^2 + \lambda}}{2}$$

$\textcircled{3}$ There do not exist two different singular values $\mu_i > \mu_j (i < j)$ such that

$$\hat{\tau}_i = \frac{\mu_i - \sqrt{\mu_i^2 + \lambda}}{2}, \quad \hat{\tau}_j = \frac{\mu_j - \sqrt{\mu_j^2 + \lambda}}{2}$$

Proof

1) If statement $\textcircled{1}$ is not true, then there exist at least two components $\hat{\tau}_i, \hat{\tau}_j (i < j)$ of $\hat{\tau}$, such that $\hat{\tau}_i < \hat{\tau}_j$. We then exchange the values of these two components to obtain a new admissible choice $\tilde{\tau}$. It is easy to see that $g(\tilde{\tau}) < g(\hat{\tau})$, which leads to a contradiction. The desired result then follows.

2) If statement $\textcircled{2}$ is not true, then there exists some $j_0 > i$ such that $\mu_{j_0}^2 + \lambda > 0$ and

$$\hat{\tau}_{j_0} = \frac{\mu_{j_0} + \sqrt{\mu_{j_0}^2 + \lambda}}{2}$$

Thus the above statement gives

$$\hat{\tau}_i \geq \hat{\tau}_{j_0}$$

i. e. ,

$$\mu_i - \mu_{j_0} \geq \sqrt{\mu_i^2 + \lambda} + \sqrt{\mu_{j_0}^2 + \lambda}$$

or equivalently,

$$-(\lambda + \mu_i \mu_{j_0}) \geq \sqrt{\mu_i^2 + \lambda} \sqrt{\mu_{j_0}^2 + \lambda}$$

That means $\lambda + \mu_i \mu_{j_0} \leq 0$, which together with $\mu_{j_0} \leq \mu_i$ yields

$$\lambda + \mu_{j_0}^2 \leq \lambda + \mu_i \mu_{j_0} \leq 0$$

We then obtain a contradiction. This completes the proof of statement $\textcircled{2}$ of lemma 4.

3) If statement $\textcircled{3}$ is not true, then there exist two different singular values $\mu_{i_0} > \mu_{j_0} (i_0 < j_0)$ such that

$$\hat{\tau}_{i_0} = \frac{\mu_{i_0} - \sqrt{\mu_{i_0}^2 + \lambda}}{2}, \quad \hat{\tau}_{j_0} = \frac{\mu_{j_0} - \sqrt{\mu_{j_0}^2 + \lambda}}{2}$$

Thus it follows from statement $\textcircled{1}$ just proved that $\hat{\tau}_{i_0} \geq \hat{\tau}_{j_0}$, which results in $\lambda \geq 0$ by some simple calculation. When $\lambda = 0$, $\hat{\tau}_{i_0} = 0$ which contradicts the admissible condition (10). When $\lambda > 0$, it is clear that $\hat{\tau}_{i_0} < 0$, $\hat{\tau}_{j_0} < 0$. We then change the components of $\hat{\tau}$ at the positions i_0 and j_0 into $-\hat{\tau}_{i_0}$, $-\hat{\tau}_{j_0}$, respectively, and that gives us a new admissible choice $\tilde{\tau}$. It is evident that $g(\tilde{\tau}) < g(\hat{\tau})$. This is a contradiction. The desired result is therefore verified.

By virtue of lemmas 1 to 4, we can obtain the following result directly.

Theorem 1 Let the nonsingular matrix \mathbf{A} admit the singular value decomposition (4) and assume that its smallest singular value has multiplicity n_0 , i. e. , $\mu_{n-n_0+1} = \dots = \mu_n$. Then there exists some nonnegative integer $k (0 \leq k \leq n_0)$ and some real constant λ satisfying the equation

$$\prod_{i=1}^{n-k} \frac{\mu_i + \sqrt{\mu_i^2 + \lambda}}{2} \left(\frac{\mu_n - \sqrt{\mu_n^2 + \lambda}}{2} \right)^k = \text{sgn}(\det(\mathbf{A})) \quad (13)$$

such that the following matrix solves problem (1):

$$\hat{\mathbf{X}} = \mathbf{P}_1 \mathbf{D} \mathbf{P}_2$$

with

$$\mathbf{D} = \text{diag} \left(\frac{\mu_1 + \sqrt{\mu_1^2 + \lambda}}{2}, \dots, \frac{\mu_{n-k} + \sqrt{\mu_{n-k}^2 + \lambda}}{2}, \frac{\mu_n - \sqrt{\mu_n^2 + \lambda}}{2}, \dots, \frac{\mu_n - \sqrt{\mu_n^2 + \lambda}}{2} \right)$$

Proof According to lemma 1 and lemma 2, problem (1) can be converted into problem (11), the solution of which should take the form (12) by lemma 3. Furthermore, by lemma 4 only the following combinations may solve problem (11):

$$\tau_1 = \frac{\mu_1 + \sqrt{\mu_1^2 + \lambda}}{2}, \dots, \tau_{n-l} = \frac{\mu_{n-l} + \sqrt{\mu_{n-l}^2 + \lambda}}{2},$$

$$\tau_{n-l+1} = \frac{\mu_n - \sqrt{\mu_n^2 + \lambda}}{2}, \dots, \tau_n = \frac{\mu_n - \sqrt{\mu_n^2 + \lambda}}{2}$$

where $0 \leq l \leq n_0$. We then obtain the result of theorem 1 easily.

3 Numerical Results

We first transfer Eq. (13) into

$$f(\lambda) \equiv \prod_{i=1}^{n-k} \frac{\mu_i + \sqrt{\mu_i^2 + \lambda}}{2} \left(\frac{\sqrt{\mu_n^2 + \lambda} - \mu_n}{2} \right)^k - (-1)^k \text{sgn}(\det(\mathbf{A})) = 0 \quad (14)$$

It is clear that $f(\lambda)$ is strictly increasing as $\lambda \in (0, +\infty)$. Hence, only if $f(0) < 0$, Eq. (14) has one (only one) solution in $(0, +\infty)$.

According to theorem 1, we then have the following algorithm for solving problem (1):

① Compute the singular value decomposition of A to obtain P_1, P_2 and Λ used in Eq. (4).

② Let n_0 be the multiplicity of the smallest singular value of A . For $k = 0, 1, \dots, n_0$, compute the whole set of solutions of Eq. (14) in the interval $[-\mu_n^2, 0]$; if $f(0) < 0$, then also compute the unique solution of Eq. (14) in $(0, +\infty)$.

③ For each λ obtained above, compute the related τ by the formula given in theorem 1, and then compute $g(\tau)$ by Eq. (9). Determine those λ which correspond to the minimal value of $g(\tau)$ just obtained.

④ For each of these τ , compute $\hat{X} = P_1 D P_2$, where D is a diagonal matrix with the i -th entry taking the value τ_i , the i -th component of τ . Then \hat{X} solves problem (1).

We next use the algorithm to solve a concrete problem. We randomly generate a matrix as follows:

$$A = \begin{bmatrix} 0.8462 & 0.6721 & 0.6813 \\ 0.5252 & 0.8381 & 0.3795 \\ 0.2026 & 0.0196 & 0.8318 \end{bmatrix}$$

The related singular value decomposition is

$$P_1 = \begin{bmatrix} 0.7336 & -0.0557 & -0.6773 \\ 0.5790 & -0.4706 & 0.6658 \\ 0.3558 & 0.8806 & 0.3130 \end{bmatrix}$$

$$P_2 = \begin{bmatrix} 0.5760 & 0.5693 & 0.5867 \\ -0.1724 & -0.6170 & 0.7679 \\ -0.7991 & 0.5434 & 0.2572 \end{bmatrix}$$

$$\Lambda = \text{diag}(1.7309, 0.6719, 0.2004)$$

In this case, $n_0 = 1$. By computation we also find that only as $k=0$ Eq. (14) is solvable and $\lambda = 0.9082$, $\tau = \{1.8534, 0.9190, 0.5871\}$, $g(\tau) = 0.2256$. Thus

$$\hat{X} = P_1 D P_2 = \begin{bmatrix} 1.1097 & 0.5894 & 0.6561 \\ 0.3803 & 1.0902 & 0.3984 \\ 0.0935 & -0.0240 & 1.0556 \end{bmatrix}$$

solves problem (1) with the above matrix A .

References

- [1] Warner F W. *Foundations of differentiable manifolds and Lie groups* [M]. Berlin: Springer, 1983.
- [2] Andersson L E, Elfving T. A constrained procrustes problem [J]. *SIAM J Matrix Anal Appl*, 1997, **18**: 124 – 139.
- [3] Arun K S. A unitarily constrained total least squares problem in signal processing [J]. *SIAM J Matrix Anal Appl*, 1992, **13**(3): 729 – 745.
- [4] Higham N J. Computing a nearest symmetric positive semidefinite matrix [J]. *Linear Algebra and Its Applications*, 1988, **103**: 103 – 118.
- [5] Higham N J. The symmetric procrustes problem [J]. *BIT*, 1988, **28**(1): 133 – 143.
- [6] Umeyama S. Least-squares estimation of transformation parameters between two point patterns [J]. *IEEE Trans Pattern Anal Mach Intelligence*, 1991, **13**(4): 376 – 380.
- [7] Zietak K. On approximation problems with zero-trace matrices [J]. *Linear Algebra and Its Applications*, 1996, **247**: 169 – 183.
- [8] Arnold V I. *Mathematical methods of classical mechanics* [M]. Berlin: Springer, 1978.
- [9] Cyganski D, Orr J A. Applications of tensor theory to object recognition and orientation determination [J]. *IEEE Trans Pattern Anal Mach Intelligence*, 1985, **7**(6): 663 – 673.
- [10] Golub G H, van Loan C F. *Matrix computations* [M]. Baltimore: The Johns Hopkins University Press, 1989.

求解最近保体矩阵问题的若干研究

沐建飞 黄建国

(上海交通大学数学系, 上海 200240)

摘要: 研究了求解给定矩阵的最近保体矩阵问题. 首先导出该问题解所必须满足的一个矩阵方程, 然后用奇异值分解方法求解该矩阵方程; 并获得了该问题解的其他更进一步的刻画条件. 利用这些结果建立了一个求解算法, 并通过数值算例说明了该算法的有效性.

关键词: 保体矩阵; 矩阵拟合问题; 奇异值分解

中图分类号: O241.1