



Minimum-correction second-moment matching: theory, algorithms and applications

Jing Lin¹ · Pierre F. J. Lermusiaux¹

Received: 18 September 2019 / Revised: 7 November 2020 / Accepted: 24 December 2020 /
Published online: 6 February 2021

© The Author(s), under exclusive licence to Springer-Verlag GmbH, DE part of Springer Nature 2021

Abstract

We address the problem of finding the closest matrix \tilde{U} to a given U under the constraint that a prescribed second-moment matrix \tilde{P} must be matched, i.e. $\tilde{U}^T \tilde{U} = \tilde{P}$. We obtain a closed-form formula for the unique global optimizer \tilde{U} for the full-rank case, that is related to U by an SPD (symmetric positive definite) linear transform. This result is generalized to rank-deficient cases as well as to infinite dimensions. We highlight the geometric intuition behind the theory and study the problem's rich connections to minimum congruence transform, generalized polar decomposition, optimal transport, and rank-deficient data assimilation. In the special case of $\tilde{P} = I$, minimum-correction second-moment matching reduces to the well-studied optimal orthonormalization problem. We investigate the general strategies for numerically computing the optimizer and analyze existing polar decomposition and matrix square root algorithms. We modify and stabilize two Newton iterations previously deemed unstable for computing the matrix square root, such that they can now be used to efficiently compute both the orthogonal polar factor and the SPD square root. We then verify the higher performance of the various new algorithms using benchmark cases with randomly generated matrices. Lastly, we complete two applications for the stochastic Lorenz-96 dynamical system in a chaotic regime. In reduced subspace tracking using dynamically orthogonal equations, we maintain the numerical orthonormality and continuity of time-varying base vectors. In ensemble square root filtering for data assimilation, the prior samples are transformed into posterior ones by matching the covariance given by the Kalman update while also minimizing the corrections to the prior samples.

Mathematics Subject Classification 65F25 · 15A23 · 60H15

✉ Pierre F. J. Lermusiaux
pierre@mit.edu

Jing Lin
linjing@mit.edu

¹ Department of Mechanical Engineering, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

1 Introduction

The second-moment matrix of $U \in \mathcal{M}_{m \times n}(\mathbb{R})$ ¹ defined as $P = U^T \Gamma_m U \in \mathcal{M}_{n \times n}$ contains all the pairwise inner products of U 's columns weighted by Γ_m . When we view U 's columns as vectors in an inner product space, P is called the Gram matrix. When the U 's rows are samples of a zero-mean random vector, P becomes the sample covariance matrix with $\Gamma_m = (1/m)I$. Here we unify them as second-moment matrices. In many computational problems, we encounter the task of correcting a given matrix U to some \tilde{U} that matches a prescribed second-moment matrix \tilde{P} , i.e. $\tilde{P} = \tilde{U}^T \Gamma_m \tilde{U}$. However, such a \tilde{U} is not unique. When there is no physical information that favors one choice over another, a natural approach is to aim for a minimal correction $\tilde{U} - U$ to avoid introducing numerical artifacts.

Precisely, given a target symmetric positive definite (SPD) second-moment matrix \tilde{P} , as well as an inner product on \mathbb{R}^m and on \mathbb{R}^n defined by the SPD weight matrix $\Gamma_m \in \mathcal{M}_{m \times m}$ and $\Gamma_n \in \mathcal{M}_{n \times n}$, respectively, we want to solve the optimization

$$\arg \min_{\tilde{U} \in \mathcal{M}_{m \times n}: \tilde{U}^T \Gamma_m \tilde{U} = \tilde{P}} \|\tilde{U} - U\|_{F, \Gamma_m, \Gamma_n}^2, \tag{1}$$

where $\|\cdot\|_{F, \Gamma_m, \Gamma_n}$ is the Frobenius norm weighted by Γ_m and Γ_n :

$$\begin{aligned} \|V\|_{F, \Gamma_m, \Gamma_n} &\triangleq \text{tr}(\sqrt{\Gamma_n} V^T \Gamma_m V \sqrt{\Gamma_n}^{-T}) \\ &= \text{tr}(\sqrt{\Gamma_m} V \Gamma_n V^T \sqrt{\Gamma_m}^{-T}) = \|V^T\|_{F, \Gamma_n, \Gamma_m}. \end{aligned} \tag{2}$$

Here \sqrt{A} denotes a matrix square root of an SPD matrix A , viewed as a self-adjoint positive operator, i.e. $\sqrt{A}^T \sqrt{A} = A$. Throughout this paper, ‘‘matrix square root’’ will refer to this definition (rather than $\sqrt{A^2} = A$), unless otherwise mentioned. Such a square root is not unique and is subject to the unitary freedom, i.e. if R is a square root of A , the set $\{QR : Q^T Q = I\}$ contains all square roots of A . The choice of \sqrt{A} can be arbitrary but must be consistent. Moreover, we denote by $A^{1/2}$ the unique SPD square root of A .

For the optimization problem (1), we can always eliminate the inner product’s weights Γ_m and Γ_n by variable substitution. If we take $W = \sqrt{\Gamma_m} U \sqrt{\Gamma_n}^{-T}$ and $\tilde{W} = \sqrt{\Gamma_m} \tilde{U} \sqrt{\Gamma_n}^{-T}$, then (1) reduces to

$$\arg \min_{\tilde{W} \in \mathcal{M}_{m \times n}: \tilde{W}^T \tilde{W} = \sqrt{\Gamma_n} \tilde{P} \sqrt{\Gamma_n}^{-T}} \|\tilde{W} - W\|_F^2, \tag{3}$$

where $\|V\|_F = \text{tr}(V^T V)$ is the unweighted Frobenius norm and $\sqrt{\Gamma_n} \tilde{P} \sqrt{\Gamma_n}^{-T}$ becomes the new target second-moment matrix. Therefore, without loss of generality, we will restrict ourselves to the simplified problem

¹ Since we will exclusively focus on real matrices in this paper, the explicit specification ‘‘(\mathbb{R})’’ of the field of matrix elements will be dropped hereafter.

$$\arg \min_{\tilde{U} \in \mathcal{M}_{m \times n}: \tilde{U}^T \tilde{U} = \tilde{P}} \|\tilde{U} - U\|_F^2. \quad (4)$$

Note that we can also scale \tilde{U} to make $\tilde{P} = I$, but then $\Gamma_n \neq I$ in general. We choose to have $\Gamma_n = I$ and \tilde{P} arbitrary because it renders the problem's structure symmetric and leads to more insightful interpretations as will be seen in Sect. 2. How the two ways of scaling are connected will be mentioned in Sect. 2.3. Besides \tilde{P} being SPD, we further assume that $m \geq n$ and U has full column rank, so $P = U^T U$ is SPD. The rank-deficiency complications due to the violation of these conditions are addressed in Sect. 2.4.

Minimum-correction second-moment matching has many applications. In the special case of $\tilde{P} = I$, the task reduces to orthonormalization. Many algorithms involve tracking a varying orthogonal matrix [14,45], such as in optimization [1] or solving matrix differential equations with orthogonality constraints [22, sec. IV.9]. There does exist orthogonality-preserving algorithms for some applications. For example, [56] studies an orthogonality-preserving curvilinear search algorithm for optimization on a Stiefel manifold [14]. For time integration, [11,24] investigates orthogonality-preserving Gauss–Legendre Runge–Kutta schemes. However, such algorithms are typically constrained and the choices limited. Most matrix update algorithms prioritize other goals such as steepest descent optimization [50] or accurate time integration [10,30,39,46,53]. In such cases, after an orthogonal matrix is updated, a deviation from orthogonality is usually incurred by the numerical discretization of the matrix's continuous evolution. Therefore, we need to orthonormalize the matrix in such a way that this *entry-wise* continuous evolution is preserved [18,21]. A natural idea is thus to find the closest point to the updated matrix on the Stiefel manifold [2,11,25]. We will see a test case of this kind in Sect. 5.

In the general case with $\tilde{P} \neq I$, one application is minimum-correction covariance matching. For example, in an ensemble Kalman filter (EnKF) [16] for data assimilation, given the prior samples of a random vector in the rows of U and some observation data, the goal is to obtain posterior samples \tilde{U} such that their empirical mean and covariance match those of the posterior distribution obtained by the Kalman update. A variant proposed by [41] consists of updating the mean in the same way as EnKF but making the posterior sample variation \tilde{U} (with mean removed) as close as possible to the prior counterpart U under the constraint that the empirical covariance $(1/m)\tilde{U}^T \tilde{U}$ must match the one obtained by the Kalman update. We will use such a test case in Sect. 6.

Although the minimum-correction orthonormalization is relatively well understood with efficient algorithms developed, its general second-moment matching counterpart has not yet been studied and analyzed in a unified and systematic way. We aim to fill this gap in the present paper.

In Sect. 2, we solve the optimization (4) analytically with a constructive proof (Sect. 2.1 and “Appendix A”), provide geometric intuition, solve rank deficient cases (2.4), and discuss the connections to other topics including minimum congruence transform (Sect. 2.2) and generalized polar decomposition (Sect. 2.3). We also address the generalization to infinite dimensions in “Appendix B” and connect to optimal trans-

port in ‘‘Appendix C’’. In Sect. 3, we obtain numerical strategies for computing the optimal solution and review existing algorithms for polar decomposition and matrix square root that play an essential role in minimum-correction second-moment matching. We then modify and stabilize two Newton iterations deemed unstable previously, that can now be used to efficiently compute both the orthogonal polar factor and the SPD square root. We compare algorithms in terms of accuracy, cost and robustness. Results are benchmarked using randomly generated matrices in Sect. 4. In Sect. 5, we show that these algorithms can be used in a subspace tracking method for stochastic ODEs and PDEs, the dynamically orthogonal equations [17,46], to maintain both the orthonormality and the time continuity of the base vectors. Furthermore, we show in Sect. 6 how the algorithms for the general case can be applied to an ensemble-based filter for data assimilation, as introduced in [41]. Lastly, we conclude in Sect. 7.

2 Theory of minimum-correction second-moment matching

2.1 Analytic expression of the global optimizer

The constrained optimization (4) has a compact feasible region and a smooth objective function, so a global optimizer exists and is a critical point of the Lagrangian. Hence, we employ the approach of Lagrangian multipliers as in [41]. We complete a different proof, adopt more compact matrix notations, and provide new extensive discussions. The results are summarized in the following theorem with the proof in ‘‘Appendix A’’.

Theorem 2.1 (Minimum-correction second-moment matching) *The global optimum of (4), assuming that $m \geq n$ and U has full column rank, is achieved by applying the unique $n \times n$ SPD linear transform to U that matches the second-moment matrix of the result \tilde{U} to the target \tilde{P} . More precisely,*

$$\arg \min_{\tilde{U} \in \mathcal{M}_{m \times n}: \tilde{U}^T \tilde{U} = \tilde{P}} \|\tilde{U} - U\|_F^2 = U A_*, \tag{5}$$

$$\|U A_* - U\|_F^2 = \text{tr} \left[P + \tilde{P} - 2 \left(\sqrt{\tilde{P}} P \sqrt{\tilde{P}}^T \right)^{1/2} \right] = \text{tr} \left[P + \tilde{P} - 2 \left(P \tilde{P} \right)^{1/2} \right] \tag{6}$$

² with the SPD linear transform A_* being

$$\begin{aligned} A_* &= \sqrt{\tilde{P}}^T \left(\sqrt{\tilde{P}} P \sqrt{\tilde{P}}^T \right)^{-1/2} \sqrt{\tilde{P}} \\ &= \sqrt{P}^{-1} \left(\sqrt{P} \tilde{P} \sqrt{P}^T \right)^{1/2} \sqrt{P}^{-T}. \end{aligned} \tag{7}$$

² For square A diagonalizable with nonnegative eigenvalues and EVD $A = V \Lambda V^{-1}$, we define $A^{1/2} \triangleq V \Lambda^{1/2} V^{-1}$. Since $P \tilde{P} = \sqrt{\tilde{P}}^{-1} \left(\sqrt{\tilde{P}} P \sqrt{\tilde{P}}^T \right) \sqrt{\tilde{P}}$ is similar to $\sqrt{\tilde{P}} P \sqrt{\tilde{P}}^T$, $P \tilde{P}$, though not symmetric in general, is diagonalizable with positive eigenvalues and EVD $P \tilde{P} = V \Lambda V^{-1}$. Therefore, $\text{tr}((P \tilde{P})^{1/2}) = \sum_{i=1}^n \sqrt{\lambda_i} = \text{tr}((\sqrt{\tilde{P}} P \sqrt{\tilde{P}}^T)^{1/2})$.

‘Note that when $\tilde{\mathbf{P}} = \mathbf{P}$, we have $\mathbf{A}_* = \mathbf{I}$, which correctly gives the null correction in this trivial case. Moreover, $\tilde{\mathbf{U}}_* = \mathbf{U}\mathbf{A}_*$ implies $\text{Col}(\tilde{\mathbf{U}}_*) = \text{Col}(\mathbf{U})$, which is particularly important for covariance matching because a zero-mean \mathbf{U} (i.e. $\mathbf{e}^T \mathbf{U} = \mathbf{0}$, where $\mathbf{e}^T = [1, \dots, 1]$) will yield a zero-mean $\tilde{\mathbf{U}}_*$.

Since $\mathcal{M}_{m \times n}$ can be equipped with the inner product $\langle \mathbf{U}, \mathbf{V} \rangle_F \triangleq \text{tr}(\mathbf{U}^T \mathbf{V})$ which induces the Frobenius norm $\|\cdot\|_F$, Theorem 2.1 indeed characterizes the orthogonal projection $\mathcal{P}_{\tilde{\mathbf{P}}}$ of any $\mathbf{U} \in \mathcal{M}_{m \times n}$ onto the $(mn - \frac{1}{2}n(n+1))$ dimensional sub-manifold $\mathcal{S}_{\tilde{\mathbf{P}}} = \{\tilde{\mathbf{U}} \in \mathcal{M}_{m \times n} : \tilde{\mathbf{U}}^T \tilde{\mathbf{U}} = \tilde{\mathbf{P}}\}$ as the unique $n \times n$ SPD linear transform that maps \mathbf{U} into the sub-manifold.

Moreover, if we invert the first expression in (7), we obtain

$$\mathbf{A}_*^{-1} = \sqrt{\tilde{\mathbf{P}}}^{-1} (\sqrt{\tilde{\mathbf{P}}}\mathbf{P}\sqrt{\tilde{\mathbf{P}}}^T)^{1/2} \sqrt{\tilde{\mathbf{P}}}^{-T},$$

which is nothing but the second expression in (7) with $\tilde{\mathbf{P}}$ and \mathbf{P} switched. This implies that $\tilde{\mathbf{U}}\mathbf{A}_*^{-1} = \mathbf{U}$ projects $\tilde{\mathbf{U}}$ orthogonally onto $\mathcal{S}_{\mathbf{P}}$, i.e.

$$\arg \min_{\mathbf{U} \in \mathcal{M}_{m \times n} : \mathbf{U}^T \mathbf{U} = \mathbf{P}} \|\mathbf{U} - \tilde{\mathbf{U}}\|_F^2 = \tilde{\mathbf{U}}\mathbf{A}_*^{-1} \tag{8}$$

with \mathbf{A}_* given by (7). Therefore, the orthogonal projection $\mathcal{P}_{\mathbf{P}}$ and $\mathcal{P}_{\tilde{\mathbf{P}}}$ specify a bijection between $\mathcal{S}_{\mathbf{P}}$ and $\mathcal{S}_{\tilde{\mathbf{P}}}$ such that

$$\mathcal{P}_{\mathbf{P}} \circ \mathcal{P}_{\tilde{\mathbf{P}}}|_{\mathcal{S}_{\tilde{\mathbf{P}}}} = \text{Id}|_{\mathcal{S}_{\tilde{\mathbf{P}}}}, \quad \mathcal{P}_{\tilde{\mathbf{P}}} \circ \mathcal{P}_{\mathbf{P}}|_{\mathcal{S}_{\mathbf{P}}} = \text{Id}|_{\mathcal{S}_{\mathbf{P}}}. \tag{9}$$

In addition, since \mathbf{A}_* depends on \mathbf{P} and $\tilde{\mathbf{P}}$ only, both $\mathcal{P}_{\tilde{\mathbf{P}}}|_{\mathcal{S}_{\tilde{\mathbf{P}}}}$ and $\mathcal{P}_{\mathbf{P}}|_{\mathcal{S}_{\mathbf{P}}}$ are characterized by a constant SPD linear transform on \mathbb{R}^n .

Furthermore, the simple case of $n = 1$ helps us build a geometric intuition of the above results, as illustrated in Fig. 1. When $n = 1$, \mathbf{P} and $\tilde{\mathbf{P}}$ reduce to a positive real number, while $\mathcal{S}_{\mathbf{P}}$ and $\mathcal{S}_{\tilde{\mathbf{P}}}$ reduce to two concentric $(m - 1)$ -spheres (or two $(m - 1)$ -ellipsoids with the same eccentricity if $\mathbf{\Gamma}_m \neq \lambda \mathbf{I}$). Hence, it is straightforward to see that the rays from the origin provide a bijection between the two spheres. This mapping is equivalent to scaling one sphere by a constant ratio. As Fig. 1 shows, the points on the two manifolds can be paired up such that in each point pair, the blue point is the closet point on the blue manifold to the red one and vice versa.

However, the case of $n = 1$ also over-simplifies some aspects of the problem and could be misleading. For example, Fig. 1 could create the illusion that the projection operator is transitive, i.e. projecting from $\mathcal{S}_{\mathbf{P}_1}$ to $\mathcal{S}_{\mathbf{P}_3}$ is equivalent to first projecting from $\mathcal{S}_{\mathbf{P}_1}$ to some intermediate $\mathcal{S}_{\mathbf{P}_2}$ and then from $\mathcal{S}_{\mathbf{P}_2}$ to $\mathcal{S}_{\mathbf{P}_3}$. We thus emphasize that in general,

$$\mathcal{P}_{\mathbf{P}_3} \circ \mathcal{P}_{\mathbf{P}_2}|_{\mathcal{S}_{\mathbf{P}_1}} \neq \mathcal{P}_{\mathbf{P}_3}|_{\mathcal{S}_{\mathbf{P}_1}}, \quad \mathbf{A}_*^{1 \rightarrow 2} \mathbf{A}_*^{2 \rightarrow 3} \neq \mathbf{A}_*^{1 \rightarrow 3}, \tag{10}$$

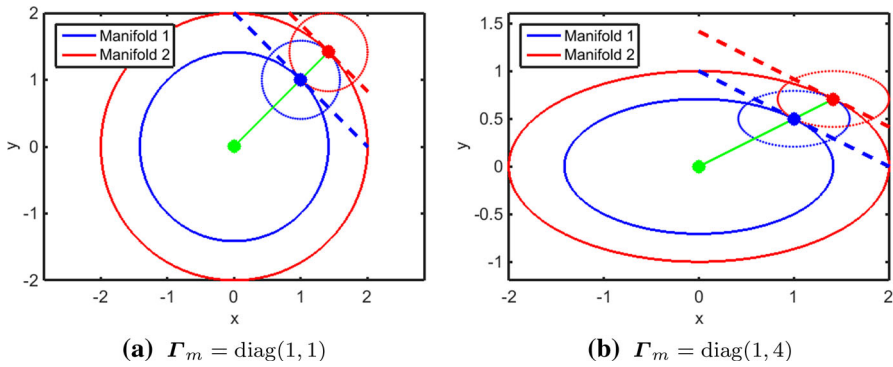


Fig. 1 A geometric picture of Theorem 2.1 for $m = 2$ and $n = 1$

since $A_*^{1 \rightarrow 3}$ is symmetric but $A_*^{1 \rightarrow 2} A_*^{2 \rightarrow 3}$ is not in general. This is especially clear when $P_2 = I$, which implies $A_*^{1 \rightarrow 2} = P_1^{-1/2}$ and $A_*^{2 \rightarrow 3} = P_3^{1/2}$ according to (7). Hence, whenever P_1 and P_3 are not commutative, $A_*^{1 \rightarrow 2} A_*^{2 \rightarrow 3}$ is not symmetric.

Finally, with (3), we obtain the global optimizer for the general setting (1).

Corollary 2.1 (General setting with inner product weighted by Γ_m and Γ_n)

$$\arg \min_{\tilde{U} \in \mathcal{M}_{m \times n}: \tilde{U}^T \Gamma_m \tilde{U} = \tilde{P}} \|\tilde{U} - U\|_{F, \Gamma_m, \Gamma_n}^2 = U(\sqrt{\Gamma_n}^{-T} A_* \sqrt{\Gamma_n}^{-T}) \tag{11}$$

$$\|\tilde{U}_* - U\|_F^2 = \text{tr}[P_\Gamma + \tilde{P}_\Gamma - 2(P_\Gamma \tilde{P}_\Gamma)^{1/2}] \tag{12}$$

$$\begin{aligned} A_* &= \sqrt{\tilde{P}_\Gamma}^{-T} (\sqrt{\tilde{P}_\Gamma} P_\Gamma \sqrt{\tilde{P}_\Gamma})^{-1/2} \sqrt{\tilde{P}_\Gamma} \\ &= \sqrt{P_\Gamma}^{-1} (\sqrt{P_\Gamma} \tilde{P}_\Gamma \sqrt{P_\Gamma}^T)^{1/2} \sqrt{P_\Gamma}^{-T}. \end{aligned} \tag{13}$$

$$\tilde{P}_\Gamma = \sqrt{\Gamma_n} \tilde{P} \sqrt{\Gamma_n}^T, \quad P_\Gamma = \sqrt{\Gamma_n} P \sqrt{\Gamma_n}^T. \tag{14}$$

We have shown previously that the choice of square root $\sqrt{\tilde{P}_\Gamma}$ and $\sqrt{P_\Gamma}$ in (13) is irrelevant. Here we can show that the choice of square root $\sqrt{\Gamma_n}$ is also irrelevant. Under the unitary freedom $\sqrt{\Gamma_n} \rightarrow Q\sqrt{\Gamma_n}$ and thus $P_\Gamma \rightarrow QP_\Gamma Q^T$, we can choose $\sqrt{QP_\Gamma Q^T}$ to be $Q\sqrt{P_\Gamma} Q^T$. Therefore, the Q 's will pair up with the Q^T 's and cancel each other, which renders the linear transform $\sqrt{\Gamma_n}^{-T} A_* \sqrt{\Gamma_n}^{-T}$ invariant. However, the choice of the inner product weights Γ_m and Γ_n themselves does affect the optimum, although they play quite different roles.

Γ_m defines the inner product on \mathbb{R}^m for computing the second-moment matrices. This inner product usually has physical meaning and is commonly determined by the underlying problem. For example, if each column of U is a discrete representation of a function $f(x)$ over some domain Ω , then $u^T \Gamma_m v$ could approximate the inner product $\int_\Omega u(x)v(x)dx$. If each row of U is a sample of a zero-mean random vector, then $\Gamma_m = \frac{1}{m} \text{diag}(1, \dots, 1)$ will make $u^T \Gamma_m v$ the empirical covariance matrix. Note

that A_* does not depend on Γ_m explicitly, but only through how \tilde{P} and P are computed based on Γ_m .

In contrast, Γ_n is more of a tunable numerical parameter. It controls how we want to weight the columns of U differently so that when we evaluate the closeness between \tilde{U} and U , we give priority to some columns over others. Such differentiation could be advantageous if, for example, some of the columns of U contain more reliable information (e.g. due to smaller numerical errors) than others, or the columns are different by several orders of magnitude and the small columns barely contribute to the closeness metric. If the magnitude is of concern, a natural choice of a scaling matrix is either $\Gamma_n = P^{-1}$ or $\Gamma_n = \tilde{P}^{-1}$:

$$\begin{aligned} \Gamma_n = P^{-1} &\implies \sqrt{\Gamma_n}^T A_* \sqrt{\Gamma_n}^{-T} = \sqrt{P}^{-1} (\sqrt{P}^{-T} \tilde{P} \sqrt{P}^{-1})^{1/2} \sqrt{P} \\ \Gamma_n = \tilde{P}^{-1} &\implies \sqrt{\Gamma_n}^T A_* \sqrt{\Gamma_n}^{-T} = \sqrt{\tilde{P}}^{-1} (\sqrt{\tilde{P}}^{-T} P \sqrt{\tilde{P}}^{-1})^{-1/2} \sqrt{\tilde{P}}. \end{aligned}$$

Surprisingly, as [41] points out, these two choices of Γ_n lead to the same optimal mapping. By the properties (8) and (9) of the projection operators, we can quickly show the equivalence between these two expressions by noticing that either one can be obtained by inverting the other and switching P with \tilde{P} . Last but not least, note that the optimal mapping $\sqrt{\Gamma_n}^T A_* \sqrt{\Gamma_n}^{-T}$ explicitly depends on Γ_n and when Γ_n is not a scalar matrix, the optimal mapping is in general not symmetric.

2.2 Minimum congruence transform

In the optimization problem (4), the candidate \tilde{U} can come from all of $\mathcal{M}_{m \times n}$, but the global optimizer $U A_*$ turns out to share the column space of U . This implies that if we restrict the candidate set to the n^2 -dimensional subspace $\{\tilde{U} \in \mathcal{M}_{m \times n} : \text{Im}(\tilde{U}) \subset \text{Im}(U)\} = \{\tilde{U} = UA : A \in \mathcal{M}_{n \times n}\}$ of the mn -dimensional $\mathcal{M}_{m \times n}$, the global optimizer is still given by (7). This is equivalent to requiring the second-moment matching to be achieved by an n -by- n linear transform. With this restriction, the objective function in (4) becomes

$$\|\tilde{U} - U\|_F^2 = \|UA - U\|_F^2 = \text{tr} \left((A - I)^T P (A - I) \right) = \|A - I\|_{F,P}^2, \tag{15}$$

where the minimum correction of U translates to the minimum operation of A , quantified by being the closest linear transform to the identity mapping. Besides, the second-moment matching constraint is now on A

$$\tilde{P} = \tilde{U}^T \tilde{U} = A^T U^T U A = A^T P A.$$

Therefore, Theorem 2.1 implies the following corollary.

Corollary 2.2 (Minimum congruence transform)

$$\arg \min_{A \in \mathcal{M}_{n \times n} : A^T P A = \tilde{P}} \|A - I\|_{F,P}^2 = A_* \tag{16}$$

with A_* given by (7) for any SPD \tilde{P} and P in $\mathcal{M}_{n \times n}$.

Note that this optimization problem is related to U only through its second-moment matrix P . We can indeed interpret it in a way that does not involve U . Since $A^T P A = \tilde{P}$ indicates that \tilde{P} is congruent with P through the congruence transform A , (16) can be viewed as finding the minimum congruence transform between two SPD matrices. Here “minimum” again refers to the minimum action of an operator.

2.3 Generalized polar decomposition

An important and familiar special case of the problem (4) is when $\tilde{P} = I$ and the task reduces to minimum-correction orthonormalization, i.e. finding the closest orthogonal matrix \tilde{U} to an arbitrary $U \in \mathcal{M}_{m \times n}$ with full column rank. In this case, the candidate solutions form the Stiefel manifold $\mathcal{S}_I = \{V \in \mathcal{M}_{m \times n} : V^T V = I\}$. Therefore, Theorem 2.1 reduces to the well-known fact that a matrix can be orthogonally projected onto the Stiefel manifold through its polar decomposition [2], summarized in the following corollary.

Corollary 2.3 (Orthogonal projection onto the Stiefel manifold by polar decomposition) *Given $U \in \mathcal{M}_{m \times n}$ with full column rank and $P = U^T U$, we have*

$$\tilde{U}_* = \arg \min_{\tilde{U} \in \mathcal{S}_I} \|\tilde{U} - U\|_F^2 = U P^{-1/2}, \tag{17}$$

where $U = \tilde{U}_* P^{1/2}$ gives the unique polar decomposition of U .

Note that Corollary 2.1 implies that the above result readily generalizes to the case with a Frobenius norm weighted by Γ_m and Γ_n ; the polar decomposition then becomes the weighted polar decomposition investigated in [57]. However, our results with $\tilde{P} \neq I$ in Theorem 2.1 indeed generalize the polar decomposition even further.

Corollary 2.4 *Given $U \in \mathcal{M}_{m \times n}$ with full column rank and $\tilde{P} \in \mathcal{M}_{n \times n}$ that is SPD, we have a unique polar-like decomposition*

$$U = \tilde{U} \Lambda \tag{18}$$

such that $\tilde{U} \in \mathcal{M}_{m \times n}$, $\tilde{U}^T \tilde{U} = \tilde{P}$ and $\Lambda \in \mathcal{M}_{n \times n}$ is SPD. Moreover, this decomposition orthogonally projects U onto \tilde{U} in the sub-manifold $\mathcal{S}_{\tilde{P}}$.

This corollary enriches the polar decomposition with the symmetry between U and \tilde{U} . Note that this symmetry is preserved by our choice of scaling that simplifies (1)–(4). Should we scale \tilde{P} rather than Γ_n to I , we would instead obtain the weighted polar decomposition in [57].

2.4 Complications due to degeneracy

In Sect. 2.1, we have restricted ourselves to the nice full-rank case with $m \geq n = \text{rank}(\tilde{P}) = \text{rank}(U)$. In practice, sometimes this is not satisfied, such as in some

scenarios of data assimilation where the number m of samples of a random vector is smaller than the dimension n of the random vector itself (see Sect. 6). Hence, here we will reveal the implications of degeneracy due to any of these full-rank assumptions being violated. Before moving on, we introduce the notations $r = \text{rank}(U)$ and $\tilde{r} = \text{rank}(\tilde{P})$.

Step 1 First of all, since $\tilde{r} \leq n$, when $m < n$, it is possible that $m < \tilde{r}$, in which case $\text{rank}(\tilde{U}^T \tilde{U}) \leq m < \tilde{r}$, so there exists no $\tilde{U} \in \mathcal{M}_{m \times n}$ whose second-moment matrix can match \tilde{P} due to rank deficiency. If $m < \tilde{r}$, an easiest fix is to find a best rank- m proxy to \tilde{P} . For example, the leading rank- m truncation of the SVD of \tilde{P} (which is also the EVD here since \tilde{P} is semi-SPD) serves as the closest rank- m approximation under any unitarily invariant norm [28]. See [28, sec. 6] for a review on other variants of a nearest lower-rank approximation. With \tilde{P} replaced by such a proxy when $m < \tilde{r}$, we can always ensure $m \geq \tilde{r}$ for the modified problem.

Step 2 Next, given $m \geq \tilde{r}$, if $\tilde{r} < n$, a variable substitution by a projection from \mathbb{R}^n onto $\text{Row}(\tilde{P})$ will eliminate the rank deficiency in \tilde{P} . If we construct \tilde{Z} with its columns forming an orthonormal basis of $\text{Row}(\tilde{P})$, we have $\tilde{P} = \tilde{Z} \tilde{P}_{\tilde{r}} \tilde{Z}^T$ such that $\tilde{P}_{\tilde{r}} = \tilde{Z}^T \tilde{P} \tilde{Z} \in \mathcal{M}_{\tilde{r} \times \tilde{r}}$ has full rank. To have $\tilde{U}^T \tilde{U} = \tilde{P}$, we must require $\text{Row}(\tilde{U}) = \text{Row}(\tilde{P})$, so \tilde{U} can be uniquely represented by $\tilde{S} \in \mathcal{M}_{m \times \tilde{r}}$ such that $\tilde{U} = \tilde{S} \tilde{Z}^T$. If we complete \tilde{Z} into an orthogonal $[\tilde{Z}, \tilde{Z}_\perp] \in \mathcal{M}_{n \times n}$, since

$$\|\tilde{U} - U\|_F^2 = \|(\tilde{U} - U)[\tilde{Z}, \tilde{Z}_\perp]\|_F^2 = \|\tilde{S} - U \tilde{Z}\|_F^2 + \|U \tilde{Z}_\perp\|_F^2,$$

the original optimization problem (4) can be equivalently formulated as

$$\arg \min_{\tilde{S} \in \mathcal{M}_{m \times \tilde{r}}: \tilde{S}^T \tilde{S} = \tilde{Z}^T \tilde{P} \tilde{Z}} \|\tilde{S} - S\|_F^2 \tag{19}$$

with $S = U \tilde{Z}$. Since we already have $m \geq \tilde{r}$, with this further manipulation, we can always reduce a degenerate problem into one with $m \geq n = \text{rank}(\tilde{P})$.

Step 3 Finally, given $m \geq n = \text{rank}(\tilde{P})$, let's consider the case of $\text{rank}(U) = r < n$. If we construct Z with its columns forming an orthonormal basis of $\text{Row}(U)$, U can be uniquely represented by a full-column-rank $W \in \mathcal{M}_{m \times r}$ such that $U = W Z^T$. Moreover, since Z can be augmented to an orthogonal $[Z, Z_\perp] \in \mathcal{M}_{n \times n}$, \tilde{U} can be uniquely represented by $\tilde{W} = \tilde{U} Z$ and $\tilde{W}_\perp = \tilde{U} Z_\perp$, which yield $\tilde{U} = \tilde{W} Z^T + \tilde{W}_\perp Z_\perp^T$ and

$$\|\tilde{U} - U\|_F^2 = \|(\tilde{U} - U)[Z, Z_\perp]\|_F^2 = \|\tilde{W} - W\|_F^2 + \|\tilde{W}_\perp\|_F^2.$$

Note that $\|\tilde{W}_\perp\|_F^2 = \text{tr}(Z_\perp^T \tilde{P} Z_\perp)$ is a constant thanks to the second-moment constraint. Therefore, the original problem (4) can be reduced to

$$\arg \min_{\tilde{W}, \tilde{W}_\perp: \tilde{U}^T \tilde{U} = \tilde{P}} \|\tilde{W} - W\|_F^2, \tag{20}$$

where $\tilde{U}^T \tilde{U} = \tilde{P}$ if and only if

$$\tilde{W}^T \tilde{W} = Z^T \tilde{P} Z, \quad \tilde{W}^T \tilde{W}_\perp = Z^T \tilde{P} Z_\perp, \quad \tilde{W}_\perp^T \tilde{W}_\perp = Z_\perp^T \tilde{P} Z_\perp. \tag{21}$$

Note the subtlety that although the objective function depends on \tilde{W} only, it is not obvious whether we can decouple the optimization over \tilde{W} and that over \tilde{W}_\perp , because it is possible that for the optimizer \tilde{W}_* to

$$\arg \min_{\tilde{W} \in \mathcal{M}_{m \times r}: \tilde{W}^T \tilde{W} = Z^T \tilde{P} Z} \|\tilde{W} - W\|_F^2, \tag{22}$$

there exists no \tilde{W}_\perp such that the second and third constraint in (21) are satisfied. These two matrix constraints consist of $r(n-r)$ linear scalar equations and $(n-r)(n-r+1)/2$ quadratic ones, which amount to $(n-r)(n+r+1)/2$ independent scalar constraints. On the other hand, we have $(n-r)m$ degrees of freedom in \tilde{W}_\perp . Therefore, under our assumption of $m \geq n \geq (r+1)$, we have $(n-r)(n+r+1)/2 \leq (n-r)n \leq (n-r)m$ so there should be enough degrees of freedom to ensure the existence of a valid \tilde{W}_\perp if all the constraints are compatible with each other. In the following, we will show that indeed for any \tilde{W} satisfying the first constraint in (21), there exists a valid \tilde{W}_\perp satisfying the other two. Moreover, we will characterize the set of all feasible \tilde{W}_\perp 's. This justifies reducing (20)–(22), which is of exactly the same form as (4) and does satisfy all the full-rank assumptions we have made for (4).

To identify all the solutions satisfying the last two constraints in (21), first we introduce the notation

$$\begin{bmatrix} B_{11} & B_{12} \\ B_{12}^T & B_{22} \end{bmatrix} \triangleq \begin{bmatrix} \tilde{W}^T \\ \tilde{W}_\perp^T \end{bmatrix} [\tilde{W}, \tilde{W}_\perp] = \begin{bmatrix} Z^T \\ Z_\perp^T \end{bmatrix} \tilde{P} [Z, Z_\perp].$$

The general solution to the linear constraint $\tilde{W}^T \tilde{W}_\perp = B_{12}$ can be written as $\tilde{W}_\perp = (\tilde{W}^T)^+ B_{12} + \tilde{W}_c V$ with “+” denoting the Moore-Penrose pseudo-inverse (a.k.a. the generalized inverse), the columns of $\tilde{W}_c \in \mathcal{M}_{m \times (m-r)}$ forming an orthonormal basis of $\text{Col}(\tilde{W})^\perp$, and $V \in \mathcal{M}_{(m-r) \times (n-r)}$ arbitrary. Since B_{11} has full rank r , \tilde{W} has full column rank r and thus $\tilde{W}^+ = (\tilde{W}^T \tilde{W})^{-1} \tilde{W}^T = B_{11}^{-1} \tilde{W}^T$ and $(\tilde{W}^T)^+ = (\tilde{W}^+)^T = \tilde{W} B_{11}^{-1}$. Hence, $\tilde{W}_\perp = \tilde{W} B_{11}^{-1} B_{12} + \tilde{W}_c V$. Now the quadratic constraint $\tilde{W}_\perp^T \tilde{W}_\perp = B_{22}$ is equivalent to $V^T V = B_{22} - B_{12}^T B_{11}^{-1} B_{12}$, where the right hand side is an SPD Schur complement. A valid V exists if and only if $(m-r) \geq (n-r) \geq 1$. Therefore, the above analysis confirms our previous speculation based on comparing the number of degrees of freedom with the number of constraints. Besides, when a valid V exists, it is not unique.

We summarize the above three steps in the following theorem.

Theorem 2.2 (Degenerate counterpart of Theorem 2.1) *Given $U \in \mathcal{M}_{m \times n}$ with $\text{rank}(U) = r$ and $\tilde{P} \in \mathcal{M}_{n \times n}$ with $\text{rank}(\tilde{P}) = \tilde{r} \leq m$, we have*

$$\arg \min_{\tilde{U} \in \mathcal{M}_{m \times n}: \tilde{U}^T \tilde{U} = \tilde{P}} \|\tilde{U} - U\|_F^2 = \tilde{U}_* = \tilde{S}_* \tilde{Z}^T = (\tilde{W}_* Z^T + \tilde{W}_{\perp*} Z_{\perp}^T) \tilde{Z}^T, \tag{23}$$

where $\tilde{S}_* = \tilde{W}_* Z^T + \tilde{W}_{\perp*} Z_{\perp}^T \in \mathcal{M}_{m \times \tilde{r}}$. Moreover, the columns of $\tilde{Z} \in \mathcal{M}_{n \times \tilde{r}}$ form an orthonormal basis of $\text{Row}(\tilde{P})$ and the columns of $Z \in \mathcal{M}_{\tilde{r} \times r'}$ form an orthonormal basis of $\text{Row}(U \tilde{Z})$ with

$$r' = \tilde{r} - \dim(\text{Row}(\tilde{P}) \cap \text{Row}(U)^\perp) = r - \dim(\text{Row}(U) \cap \text{Row}(\tilde{P})^\perp). \tag{24}$$

Finally, \tilde{W}_* and $\tilde{W}_{\perp*}$ are given by

$$\tilde{W}_* = \arg \min_{\tilde{W} \in \mathcal{M}_{m \times r'}: \tilde{W}^T \tilde{W} = Z^T \tilde{Z}^T \tilde{P} \tilde{Z} Z} \|\tilde{W} - U \tilde{Z} Z\|_F^2 \tag{25}$$

and $\tilde{W}_{\perp*} = \tilde{W}_* B_{11}^{-1} B_{12} + \tilde{W}_{c*} V_* \in \mathcal{M}_{m \times (\tilde{r} - r')}$ with the columns of $\tilde{W}_{c*} \in \mathcal{M}_{m \times (m - r')}$ forming an orthonormal basis of $\text{Col}(\tilde{W}_*)^\perp$ and $V_* \in \mathcal{M}_{(m - r') \times (\tilde{r} - r')}$ arbitrary as long as $V_*^T V_* = B_{22} - B_{12}^T B_{11}^{-1} B_{12}$. The B_{ij} blocks are defined as:

$$\begin{bmatrix} B_{11} & B_{12} \\ B_{12}^T & B_{22} \end{bmatrix} \triangleq \begin{bmatrix} \tilde{W}^T \\ \tilde{W}_{\perp}^T \end{bmatrix} [\tilde{W}, \tilde{W}_{\perp}] = \begin{bmatrix} Z^T \\ Z_{\perp}^T \end{bmatrix} \tilde{Z}^T \tilde{P} \tilde{Z} [Z, Z_{\perp}].$$

Besides, the global minimum can be expressed in exactly the same way as in (6):

$$\|\tilde{U}_* - U\|_F^2 = \text{tr}[P + \tilde{P} - 2(P\tilde{P})^{1/2}]. \tag{26}$$

Proof Everything other than (26) ensues readily from knitting into one piece step 2 and step 3 shown right before we state this theorem. To prove (26), notice that

$$\|\tilde{U}_* - U\|_F^2 - \text{tr}(P) - \text{tr}(\tilde{P}) = -2\text{tr}(U^T \tilde{U}_*) = -2\text{tr}((U \tilde{Z})^T \tilde{S}_*) = -2\text{tr}((U \tilde{Z} Z)^T \tilde{W}_*).$$

Since Theorem 2.1 applies to (25), the last expression in the above can be computed by (6) to yield $\|\tilde{U}_* - U\|_F^2 = \text{tr}(P) + \text{tr}(\tilde{P}) - 2\text{tr}[(P_{r'} \tilde{P}_{r'})^{1/2}]$, where $P_{r'} = Z^T \tilde{Z}^T P \tilde{Z} Z$ and $\tilde{P}_{r'} = Z^T \tilde{Z}^T \tilde{P} \tilde{Z} Z$. Hence, what remains to be shown is $\text{tr}[(P_{r'} \tilde{P}_{r'})^{1/2}] = \text{tr}[(P \tilde{P})^{1/2}]$ and it suffices to show that $P_{r'} \tilde{P}_{r'}$ and $P \tilde{P}$ share the same nonzero eigenvalues with the same algebraic multiplicity. To demonstrate this, first compute

$$P_{r'} \tilde{P}_{r'} = Z^T \tilde{Z}^T P \tilde{Z} (Z Z^T \tilde{Z}^T \tilde{P}) \tilde{Z} Z = Z^T \tilde{Z}^T P (\tilde{Z} \tilde{Z}^T \tilde{P}) \tilde{Z} Z = Z^T \tilde{Z}^T P \tilde{P} \tilde{Z} Z$$

using the property of the orthogonal projections $Z Z^T$ and $\tilde{Z} \tilde{Z}^T$. Next, notice that $(Z^T \tilde{Z}^T P)(\tilde{P} \tilde{Z} Z)$ and $(\tilde{P} \tilde{Z} Z)(Z^T \tilde{Z}^T P)$ have identical nonzero eigenvalues with the same algebraic multiplicity. Since $(\tilde{P} \tilde{Z} Z)(Z^T \tilde{Z}^T P) = \tilde{P} \tilde{Z} \tilde{Z}^T P = \tilde{P} P$, we conclude that the nonzero spectrum of $P_{r'} \tilde{P}_{r'}$ and of $P \tilde{P}$ coincide in terms of eigenvalues and their algebraic multiplicity, which completes the proof. \square

Remark 1 \tilde{U}_* is unique if and only if $r' = \tilde{r}$, i.e. $\text{Row}(\tilde{P}) \cap \text{Row}(U)^\perp = 0$. In this case, step 3 is not needed. Given $r' = \tilde{r}$, the reverse problem seeking the closest point to \tilde{U}_* with second-moment matrix P will also have a unique solution if and only if $r' = r$, i.e. $\text{Row}(U) \cap \text{Row}(\tilde{P})^\perp = 0$. Therefore, the elements in $\mathcal{S}_P = \{U \in \mathcal{M}_{m \times n} : U^T U = P\}$ and those in $\mathcal{S}_{\tilde{P}}$ can be paired up bijectively and unambiguously by this closest-point correspondence if and only if $\text{rank}(P) = \text{rank}(\tilde{P}) = \text{rank}(P\tilde{P})$. In practice, one common special case of such is when $\text{Row}(U) = \text{Row}(\tilde{P})$, as will be seen in Sect. 6. If this is true, besides fixing rank deficiency, step 2 also preserves the value of the objective function, i.e. $\|\tilde{S} - S\|_F^2 = \|\tilde{U}\tilde{Z} - UZ\|_F^2 \equiv \|\tilde{U} - U\|_F^2$.

Remark 2 One extreme case of rank deficiency is when $\text{Row}(\tilde{P}) \perp \text{Row}(U)$ and thus $r' = 0$. This renders the optimization (23) trivial because the objective function is now constant: $\|\tilde{U} - U\|_F^2 \equiv \text{tr}(P + \tilde{P})$, so all feasible \tilde{U} 's are equally good.

Remark 3 Although the degeneracy significantly complicates the characterization of the optimizer \tilde{U}_* , the minimum of the objective function (26) surprisingly shares the same simple formula as (6) in the non-degenerate case.

3 Algorithms for minimum-correction second-moment matching

3.1 General computation strategies

To compute \tilde{U}_* in Theorem 2.1, there are two routes. We can either compute $\tilde{U}_* = UA_*$ using

$$A_* = \sqrt{\tilde{P}}^T (\sqrt{\tilde{P}}P\sqrt{\tilde{P}})^{-1/2} \sqrt{\tilde{P}} = \sqrt{P}^{-1} (\sqrt{P}\tilde{P}\sqrt{P})^{1/2} \sqrt{P}^{-T} \tag{27}$$

or directly obtain \tilde{U}_* without explicitly forming P or computing A_* .

Computing \tilde{U}_* through A_* To take the A_* route, we need a pre-processing step $P = U^T U$, a post-processing step $\tilde{U}_* = UA_*$ and the key step (27). The drawback is that forming $P = U^T U$ indeed squares the condition number, i.e. $\kappa(P) = \kappa(U)^2$. Therefore, if U is ill-conditioned, the non- A_* route is preferred.

Computing \tilde{U}_* without A_* The non- A_* route is based on directly computing the $\mathcal{S}_{\tilde{P}}$ polar factor \tilde{U}_* of the generalized PD (polar decomposition) $U = \tilde{U}_*C$ introduced in Corollary 2.4. First, we need to reduce the rectangular PD to a square one by identifying an orthonormal basis of U 's columns, i.e. $U = VR$ with $V^T V = I$ and $R \in \mathcal{M}_{n \times n}$. This can be achieved by a Householder QR factorization [51, ch. 10]. Next we reduce computing the $\mathcal{S}_{\tilde{P}}$ factor \tilde{R}_* of the generalized PD $R = \tilde{R}_*C$ to computing the orthogonal polar factor Q of the PD $R\sqrt{\tilde{P}}^T = QS = (\tilde{R}_*\sqrt{\tilde{P}}^{-1})(\sqrt{\tilde{P}}C\sqrt{\tilde{P}}^T)$. Finally, we can assemble the desired \tilde{U}_* as $\tilde{U}_* = V\tilde{R}_* = VQ\sqrt{\tilde{P}}$.

To summarize, the A_* route is eventually reduced to computing an SPD square root, while the non- A_* route is reduced to performing a polar decomposition. These two tasks are indeed intimately related and an algorithm for one task usually corresponds to one for the other, as will be seen in the next subsection.

In the above, we have assumed the problem is non-degenerate, as required by Theorem 2.1. If that is not the case, the constructive proof of Theorem 2.2 already provides an algorithm to reduce a degenerate case to a non-degenerate one.

3.2 Algorithms for the SPD square root and polar decomposition

Given an invertible $\mathbf{R} \in \mathcal{M}_{n \times n}$ with polar decomposition $\mathbf{R} = \mathbf{Q}\mathbf{C}$, where \mathbf{Q} and \mathbf{C} are the orthogonal and SPD polar factor, respectively, then \mathbf{C} is the SPD square root of $\mathbf{R}^T\mathbf{R}$, i.e. $\mathbf{C} = (\mathbf{R}^T\mathbf{R})^{1/2}$. On the other hand, given an SPD $\mathbf{P} \in \mathcal{M}_{n \times n}$ with an arbitrary square root $\sqrt{\mathbf{P}}$, $\mathbf{P}^{1/2}$ will also be the SPD polar factor of $\sqrt{\mathbf{P}}$. Therefore, computing an SPD square root and performing a polar decomposition can be reduced to each other in exactly the same vein as how the Cholesky decomposition is related to the QR factorization.

3.2.1 Algorithms based on SVD/EVD

For the polar decomposition of $\mathbf{R} \in \mathcal{M}_{n \times n}$, a straightforward SVD-based algorithm stems from the fact that the SVD of \mathbf{R} can be readily manipulated to yield the polar decomposition:

$$\mathbf{R} = \mathbf{Q}\mathbf{\Sigma}\mathbf{V}^T = (\mathbf{Q}\mathbf{V}^T)(\mathbf{V}\mathbf{\Sigma}\mathbf{V}^T).$$

This SVD-based algorithm corresponds to using the EVD (eigenvalue decomposition) to obtain the SPD square root of an SPD $\mathbf{P} \in \mathcal{M}_{n \times n}$ by $\mathbf{P} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T = (\mathbf{V}\mathbf{\Lambda}^{1/2}\mathbf{V}^T)^2$.

Since we have efficient and numerically stable algorithms for SVD and EVD [20, 51], the above SVD/EVD-based algorithms are easy to implement and robust to use. Despite this, there are cases where a non-SVD/EVD-based algorithm is desirable. It might be because, for example, we have no access to an efficient SVD/EVD procedure, we want to have more control over the accuracy of the result, or we want to utilize some special features particular to a problem at hand to further boost the efficiency. A non-SVD/EVD-based algorithm typically involves explicitly carrying out a matrix fixed-point iteration, as will be discussed next.

3.2.2 Algorithms based on fixed-point iterations

For computing the SPD square root and polar decomposition, there are also non-SVD/EVD-based algorithms that take the form of a matrix fixed-point iteration $\mathbf{X}_{k+1} = \mathbf{F}(\mathbf{X}_k)$ with $\mathbf{F} : \mathcal{M}_{n \times n} \rightarrow \mathcal{M}_{n \times n}$ involving only matrix addition, multiplication, and inversion.

If we view an iteration as a discrete dynamical system³, the minimal requirements for a matrix iteration $\mathbf{X}_{k+1} = \mathbf{F}(\mathbf{X}_k)$ to be usable in practice for seeking a target matrix \mathbf{X}_* include:

1. \mathbf{X}_* is a neutrally stable fixed point of \mathbf{F} .

³ See [48, sec. 5.1] for the definition of terms related to dynamical systems.

2. A subset \mathcal{A}_{ini} of the attraction basin $\mathcal{A}_{F, X_*} \subset \mathcal{M}_{n \times n}$ can be identified.

Here by definition, the attraction basin contains all the points that will converge to X_* under F , i.e. $X \in \mathcal{A}_{F, X_*}$ if and only if $X_k \rightarrow X_*$ with $X_0 = X$ and $X_{k+1} = F(X_k)$. If these two requirements are met, for any initial guess $X_0 \in \mathcal{A}_{\text{ini}}$, we have $X_k \rightarrow X_*$ analytically and the neutral stability of X_* guarantees that any numerical errors incurred in carrying out $X_{k+1} = F(X_k)$ will stay bounded once X_k approaches X_* , i.e. entering some ε -neighborhood of X_* .

As we shall see, usually we can find \mathcal{A}_{ini} as a stable manifold containing X_* and invariant under F . Sometimes \mathcal{A}_{ini} can even extend to infinity. However, this should not be confused with global convergence and it does not even guarantee the numerical stability of the algorithm because \mathcal{A}_{ini} as a submanifold of $\mathcal{M}_{n \times n}$ is a boundary set, i.e. it has no interior point. Hence, the numerical errors $(\hat{X}_k - X_k)$, although typically small in magnitude but in general unconstrained in structure, will drive the actual iterates \hat{X}_k 's away from \mathcal{A}_{ini} and might make them fail to converge to anything close to X_* . This is why we need the neutral stability of X_* to ensure that if the iterates manage to approach X_* before the numerical errors accumulate significantly, these errors will remain bounded hereafter.

Since the dynamical system $X_{k+1} = F(X_k)$ is n^2 -dimensional, we can reshape each X_k into an $\mathbb{R}^{n \times n}$ vector and view F as an operator from $\mathbb{R}^{n \times n}$ to itself. Denote the Jacobian matrix of F at X_* by $JF \in \mathcal{M}_{n^2 \times n^2}$ ⁴. Then a necessary condition for X_* to be neutrally stable is $\rho(JF) \leq 1$, where $\rho(\cdot)$ denotes the spectral radius, while the strict inequality $\rho(JF) < 1$ is a sufficient one. However, viewing X_k as a vector typically makes the expression of F awkwardly complicated and renders cumbersome the analysis of JF 's spectrum. [43] bypasses this difficulty elegantly by keeping X_k as a matrix and taking a functional viewpoint of $X_{k+1} = F(X_k)$. From this perspective, the counterpart of JF is the Fréchet derivative DF of the operator F at X_* . Here $DF : \mathcal{M}_{n \times n} \rightarrow \mathcal{M}_{n \times n}$ is a linear operator. Since JF and DF are essentially the same linear operator, they share the same spectrum. It turns out that identifying all the eigenvalues and eigenvectors of DF by direct observation can be made possible by suitable linear transforms.

Next, we review the most common iterations for computing the i) SPD square root of an SPD matrix and ii) orthogonal polar factor of a square matrix. To do so, we analyze and compare iterations in a novel unified framework. We also obtain new insights into their properties and propose a slight modification that surprisingly remedies some iterations deemed numerically unstable previously.

(i) A unified framework for four SPD square root Iterations

Given an SPD P , the following four iterations can be used to compute the SPD square root $P^{1/2}$ and $P^{-1/2}$ (see [27,43] and [20, sec. 9.4.2]):

$$X_{k+1} = X_k + \Delta X_k, \quad X_k \Delta X_k + \Delta X_k X_k = P - X_k^2; \tag{28}$$

$$X_{k+1} = (1/2)(X_k + X_k^{-1}P); \tag{29}$$

$$X_{k+1} = (1/2)(X_k + Y_k^{-1}), \quad Y_{k+1} = (1/2)(Y_k + X_k^{-1}); \tag{30}$$

$$X_{k+1} = X_k + (1/2)X_k(I - X_k P X_k). \tag{31}$$

⁴ The explicit dependence on X_* is omitted in the notation since it should be clear from context.

Among these iterations, (28) stems from a direct application of the Newton’s method to the matrix equation $\mathbf{G}(\mathbf{X}) = \mathbf{X}^2 - \mathbf{P} = \mathbf{0}$ with the Fréchet derivative of \mathbf{G} at \mathbf{X}_k being the linear operator $\mathbf{D}\mathbf{G}_{\mathbf{X}_k}\mathbf{X} = \mathbf{X}_k\mathbf{X} + \mathbf{X}\mathbf{X}_k$. (29) is a simplification of (28) by making the assumption $\mathbf{X}_k\Delta\mathbf{X}_k = \Delta\mathbf{X}_k\mathbf{X}_k$, which is true for certain choices of \mathbf{X}_0 as we shall see. (30) is an extension of (29) by coupling (29) with a mirror iteration to compute $\mathbf{P}^{1/2}$ and $\mathbf{P}^{-1/2}$ simultaneously. Finally, (31) is quite different from the others, as it comes from mimicking the Newton iteration for the scalar equation $g(x) = 1/x^2 - \lambda = 0$ [43].

Lemma 3.1 *The converging point and a corresponding \mathcal{A}_{ini} of feasible initial values for iterations Eqs. (28)–(31) are summarized as*

$$(28), (29) : \mathbf{X}_k \rightarrow \mathbf{P}^{1/2} \quad \text{for } \mathbf{X}_0 \in \mathcal{C}_P \triangleq \{\mathbf{X} \text{ SPD} : \mathbf{X}\mathbf{P} = \mathbf{P}\mathbf{X}\}; \tag{32}$$

$$(30) : \mathbf{X}_k \rightarrow \mathbf{P}^{1/2}, \quad \mathbf{Y}_k \rightarrow \mathbf{P}^{-1/2} \quad \text{for } \mathbf{X}_0 \in \mathcal{C}_P, \quad \mathbf{Y}_0 = \mathbf{P}^{-1}\mathbf{X}_0; \tag{33}$$

$$(31) : \mathbf{X}_k \rightarrow \mathbf{P}^{-1/2} \quad \text{for } \mathbf{X}_0 \in \{\mathbf{X} \in \mathcal{C}_P : \|\mathbf{X}\|_2 < \sqrt{3/\|\mathbf{P}\|_2}\}. \tag{34}$$

Proof Observe that $\mathbf{X}_k \in \mathcal{C}_P$ (and $\mathbf{Y}_k \in \mathcal{C}_P$ for (30)) implies $\mathbf{X}_{k+1} \in \mathcal{C}_P$ (and $\mathbf{Y}_{k+1} \in \mathcal{C}_P$), so \mathcal{C}_P is an invariant submanifold under these iterations. Therefore, there exists an EVD $\mathbf{P} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$ that also diagonalizes the iterates $\mathbf{X}_k = \mathbf{Q}\mathbf{\Lambda}_k\mathbf{Q}^T$ (and $\mathbf{Y}_k = \mathbf{Q}\mathbf{\Sigma}_k\mathbf{Q}^T$) and the matrix iteration can be decoupled into n scalar ones of the same form:

$$(28), (29) : \lambda_{k+1} = (1/2)(\lambda_k + \lambda/\lambda_k) : \lambda_k \rightarrow \sqrt{\lambda} \quad \text{for } \lambda_0 > 0 \tag{35}$$

$$(30) : \lambda_{k+1} = (1/2)(\lambda_k + 1/\sigma_k), \quad \sigma_{k+1} = (1/2)(\sigma_k + 1/\lambda_k) : \tag{36}$$

$$\lambda_k \rightarrow \sqrt{\lambda_0/\sigma_0}, \quad \sigma_k \rightarrow \sqrt{\sigma_0/\lambda_0} \quad \text{for } \lambda_0, \sigma_0 > 0$$

$$(31) : \lambda_{k+1} = (1/2)\lambda_k(3 - \lambda_k^2/\lambda) : \lambda_k \rightarrow 1/\sqrt{\lambda} \quad \text{for } 0 < \lambda_0 < \sqrt{3/\lambda} \tag{37}$$

Here $\lambda > 0$ is any diagonal entry of $\mathbf{\Lambda}$, which is also an eigenvalue of \mathbf{P} . □

Note that (36) can be reduced to two decoupled iterations of the same form as (35) due to the observation that $\lambda_k/\sigma_k \equiv \lambda_0/\sigma_0$. Since (35) is nothing but the well-known Newton square-root iteration that solves $g(x) = x^2 - \lambda = 0$ while (37) is the Newton iteration for solving $g(x) = 1/x^2 - \lambda = 0$ [43], we know both scalar iterations converge quadratically. Furthermore, this implies that the matrix iterations Eqs. (28)–(31) also converge quadratically provided their initial values satisfy Eqs. (32)–(34).

Lemma 3.1 only guarantees the analytic convergence when roundoff errors are absent. To ensure the stability of an iteration $\mathbf{X}_{k+1} = \mathbf{F}(\mathbf{X}_k)$, we require the Fréchet derivative $\mathbf{D}\mathbf{F}$ of the operator \mathbf{F} at the limit \mathbf{X}_* to have spectral radius $\rho(\mathbf{D}\mathbf{F}) \leq 1$.

Lemma 3.2 *The Fréchet derivatives for Eqs. (28)–(31) are summarized as*

$$(28) : \mathbf{D}\mathbf{F}\mathbf{X} = \mathbf{0}; \tag{38}$$

$$(29), (31) : \mathbf{D}\mathbf{F}\mathbf{X} = (1/2)(\mathbf{X} - \mathbf{\Lambda}^{-1/2}\mathbf{X}\mathbf{\Lambda}^{1/2}); \tag{39}$$

$$(30) : \mathbf{D}\mathbf{F}(\mathbf{X}, \mathbf{Y}) = (1/2)(\mathbf{X} - \mathbf{\Lambda}^{1/2}\mathbf{Y}\mathbf{\Lambda}^{1/2}, \mathbf{Y} - \mathbf{\Lambda}^{-1/2}\mathbf{X}\mathbf{\Lambda}^{-1/2}). \tag{40}$$

These linear operators have eigenvectors \mathbf{Z}_{ij} 's and eigenvalues μ_{ij} 's as below:

$$(39) : \mathbf{Z}_{ij} = \mathbf{E}_{ij}, \quad \mu_{ij} = (1/2)(1 - \sqrt{\lambda_j/\lambda_i}); \tag{41}$$

$$(40) : \mathbf{Z}_{ij} = (\mathbf{E}_{ij}, 1/\sqrt{\lambda_i\lambda_j}\mathbf{E}_{ij}), \quad \mu_{ij} = 0, \\ \mathbf{Z}'_{ij} = (\mathbf{E}_{ij}, -1/\sqrt{\lambda_i\lambda_j}\mathbf{E}_{ij}), \quad \mu'_{ij} = 1. \tag{42}$$

Here $\mathbf{E}_{ij} \in \mathcal{M}_{n \times n}$ is the matrix having its (i, j) -th entry equal to 1 as the only nonzero entry.

Proof (38) is a standard result for an exact Newton iteration, which must have $\mathbf{F}(\mathbf{X}) = \mathbf{X} - (\mathbf{D}\mathbf{G}_{\mathbf{X}})^{-1}\mathbf{G}(\mathbf{X})$. Hence, $(\mathbf{D}\mathbf{F}_{\mathbf{X}_*})\mathbf{X} = \mathbf{X} - (\mathbf{D}\mathbf{G}_{\mathbf{X}_*})^{-1}(\mathbf{D}\mathbf{G}_{\mathbf{X}_*})\mathbf{X} = \mathbf{X} - \mathbf{X} \equiv \mathbf{0}$. (39) and (40) are obtained by diagonalizing $\mathbf{P} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$ with a linear variable substitution $\mathbf{X}' = \mathbf{Q}^T\mathbf{X}\mathbf{Q}$. For example, for (39), originally we have $\mathbf{Y} = \mathbf{D}\mathbf{F}\mathbf{X} = (1/2)(\mathbf{X} - \mathbf{P}^{-1/2}\mathbf{X}\mathbf{P}^{1/2})$. After the variable substitution, $\mathbf{Y} = \mathbf{D}\mathbf{F}\mathbf{X}$ becomes $\mathbf{Y}' = \mathbf{Q}^T(\mathbf{D}\mathbf{F}(\mathbf{Q}\mathbf{X}'\mathbf{Q}^T))\mathbf{Q} = (1/2)(\mathbf{X}' - \mathbf{\Lambda}^{-1/2}\mathbf{X}'\mathbf{\Lambda}^{1/2})$, which yields (39). This simplification can be justified by the fact that the spectrum of $\mathbf{D}\mathbf{F}$ is invariant under a linear variable substitution. Note that \mathbf{X}' is not diagonal in general because in stability analysis, we can no longer assume $\mathbf{X} \in \mathcal{C}_{\mathbf{P}}$.

Since each term in (39) and (40) has only one non-diagonal matrix factor, we can readily identify all the eigenvectors \mathbf{Z}_{ij} 's and eigenvalues μ_{ij} 's by observation. \square

Therefore, if we want $\rho(\mathbf{D}\mathbf{F}) = (1/2)(\sqrt{\kappa(\mathbf{P})} - 1) < 1$ for (41), we need the 2-norm condition number $\kappa(\mathbf{P})$ to satisfy $\kappa(\mathbf{P}) < 9$. This is very unsatisfactory because it puts a stringent constraint on the \mathbf{P} 's to which (29) and (31) are applicable. [43] shows that for (31), a manual symmetrification step in the end of each iteration can alleviate this stability constraint to $\kappa(\mathbf{P}) < 17 + 6\sqrt{8}$, which is only a nonessential and minor improvement. This makes sense because $\mathcal{C}_{\mathbf{P}}$ requires being commutative with \mathbf{P} besides being SPD, so symmetrification itself does not retract a perturbed iterate back to $\mathcal{C}_{\mathbf{P}}$ in general, though may bring it closer.

On the other hand, (42) always has $\rho(\mathbf{D}\mathbf{F}) = 1$. Note that (30) has the special feature that \mathbf{P} does not show up in the iteration at all and the converging point $(\mathbf{X}_*, \mathbf{Y}_*)$ can vary with the initial value $(\mathbf{X}_0, \mathbf{Y}_0)$ continuously. Indeed, the fixed points of (30) are not isolated but form an n^2 -dimensional submanifold corresponding to the n^2 degrees of freedom in $\mathbf{X}_0\mathbf{Y}_0^{-1}$ for specifying the initial values. Moreover, the n^2 eigenvectors associated with $\mu'_{ij} = 1$ in (42) span the tangent space of this fixed-point submanifold at \mathbf{X}_* . Any perturbation along these directions will neither grow nor diminish, but will simply remain constant and eventually reflect itself on a shift in the converging point, while perturbations along the eigenvectors for $\mu_{ij} = 0$ will decay quadratically.

Before moving on, we make several remarks on the practical use of Eqs. (28)–(28). First, solving the Sylvester equation for $\Delta\mathbf{X}_k$ in (28) is nontrivial and expensive [31, sec. 13.3], which renders (28) of little practical value. Second, the extra constraint $\|\mathbf{X}_0\|_2 < \sqrt{3/\|\mathbf{P}\|_2}$ for the convergence of (31) indeed poses no difficulty because such an \mathbf{X}_0 is easy to find. Note that $\|\mathbf{X}_0\|_2 < \sqrt{3/\|\mathbf{P}\|_2}$ can be replaced by a more stringent inequality $\|\mathbf{X}_0\| < \sqrt{3/\|\mathbf{P}\|}$ since $\|\mathbf{X}_0\|_2 = \rho(\mathbf{X}_0) \leq \|\mathbf{X}_0\|$ for any operator norm $\|\cdot\|$ and symmetric \mathbf{X}_0 . Hence, we can replace the 2-norm by another

operator norm (e.g. 1-norm or ∞ -norm) that can be more easily calculated. A simple choice of X_0 can be, for example, $X_0 = (1/\sqrt{\|P\|_1})I$.

(ii) Stabilizing two unstable SPD square root iterations by combining with orthogonal polar factor iterations

Next, we review two iterations for computing the orthogonal polar factor of an invertible $R \in \mathcal{M}_{n \times n}$ (see [26] and [20, sec. 9.4.3]):

$$X_{k+1} = (1/2)(X_k + X_k^{-T}); \tag{43}$$

$$X_{k+1} = X_k + (1/2)X_k(I - X_k^T X_k), \tag{44}$$

where the initial value is set to $X_0 = R$ for both iterations. With the polar decomposition $X_0 = Q_0 S_0$, Eqs. (43) and (44) is analytically equivalent to using the same iterations with initial value $X_0 = S_0$ to compute the square root of I . Note that Eqs. (43) and (44) are almost the same as Eqs. (29) and (31) with $P = I$ except for the presence of an extra transpose. Analytically this makes no difference when the initial value X_0 is SPD since all the iterates will be symmetric. However, surprisingly it turns out that this extra transpose makes a huge difference in numerical stability because it indeed lifts the stringent $\kappa(P) < 9$ stability constraint for Eqs. (29) and (31).

We will expound this significant improvement by generalizing Eqs. (43) and (44) to the following multi-purpose (SPD square root or orthogonal polar factor) iterations:

$$X_{k+1} = (1/2)(X_k + X_k^{-T} P); \tag{45}$$

$$X_{k+1} = X_k + (1/2)X_k(I - X_k^T P X_k), \tag{46}$$

where P is some SPD matrix. Equations (45) and (46) can be viewed as the counterparts of Eqs. (29) and (31) for computing the square root of P in the sense of $X^T X = P$ rather than $X^2 = P$.

Lemma 3.3 *Suppose the initial value X_0 is invertible and has polar decomposition $X_0 = Q_0 S_0$ (or $X_0^T = Q_0^T S_0$). Then we have the following convergence properties:*

$$(45) : X_k \rightarrow Q_0 P^{1/2} \quad \text{for } X_0 = Q_0 S_0, S_0 \in \mathcal{C}_P; \tag{47}$$

$$(46) : X_k \rightarrow P^{-1/2} Q_0 \quad \text{for } X_0 = S_0 Q_0, S_0 \in \mathcal{C}_P, \|X_0\|_2 < \sqrt{3/\|P\|_2}. \tag{48}$$

In other words, X_k converges to a particular (inverse) square root of P that shares the same orthogonal polar factor with X_0 . In particular, if we set $P = I$, X_k will converge to the orthogonal polar factor Q_0 of X_0 , while if we instead set $Q_0 = I$ (i.e. $X_0 \in \mathcal{C}_P$), X_k will converge to the (inverse) SPD square root of P . Hence, these two multi-purpose iterations reflect the inherent connections between matrix square root and polar decomposition and unify the two on a algorithmic level.

Using the same methodology for stability analysis as before, we obtain:

Lemma 3.4 *The Fréchet derivative at X_* for both Eqs. (45) and (46):*

$$DFX = (1/2)(X - \Lambda^{-1/2} X^T \Lambda^{1/2}), \tag{49}$$

Table 1 Recommended algorithms for computing the minimum-correction second-moment matching solution \tilde{U}_* in Theorem 2.1. \tilde{U}_* -based algorithms avoid forming $\mathbf{P} = \mathbf{U}^T \mathbf{U}$ and directly operate on \mathbf{U}

Algorithm	Description
\mathbf{A}_* -EVD	Use (27) for \mathbf{A}_* with $\sqrt{\cdot}$ computed by Cholesky factorization and $(\cdot)^{1/2}$ computed by EVD (see Sect. 3.2.1)
\mathbf{A}_* -NtSqr	Use (27) for \mathbf{A}_* with $\sqrt{\cdot}$ computed by Cholesky factorization and $(\cdot)^{1/2}$ computed by the stablized Newton iteration (45) or (46)
\tilde{U}_* -NtPD	Perform Householder QR factorization $\mathbf{U} = \mathbf{V} \mathbf{R}$ Compute the orthogonal polar factor \mathbf{Q} of $\mathbf{R} \sqrt{\tilde{\mathbf{P}}}^T$ using the stablized Newton iteration (45) or (46). Form $\tilde{U}_* = \mathbf{V} \mathbf{Q} \sqrt{\tilde{\mathbf{P}}}$

where the diagonal $\mathbf{\Lambda}$ contains the eigenvalues of \mathbf{P} . Again, by observation, we obtain all the eigenvalues and eigenvectors for the linear operator $D\mathbf{F}$ as:

$$\begin{aligned} \mathbf{Z}_{ij} &= \mathbf{E}_{ij} + \sqrt{\lambda_i/\lambda_j} \mathbf{E}_{ji}, \quad \mu_{ij} = 0, \quad i \leq j, \\ \mathbf{Z}_{ij} &= \mathbf{E}_{ij} - \sqrt{\lambda_i/\lambda_j} \mathbf{E}_{ji}, \quad \mu_{ij} = 1, \quad i > j. \end{aligned} \tag{50}$$

Therefore, instead of the stringent constraint $\kappa(\mathbf{P}) < 9$ for Eqs. (29) and (31) due to (41), now we have $\rho(D\mathbf{F}) \equiv 1$ regardless of $\kappa(\mathbf{P})$. Moreover, the fixed points are no longer isolated and they form a $n(n - 1)/2$ -dimensional submanifold corresponding to the unitary freedom encoded by \mathbf{Q}_0 . The $n(n - 1)/2$ eigenvectors associated with eigenvalue 1 span the tangent space of this fixed-point submanifold at \mathbf{X}_* and any perturbation orthogonal to this tangent space decays quadratically to 0.

Again, as discussed in the end of part *i*), the extra initial value constraint $\|\mathbf{X}_0\|_2 < \sqrt{3}/\|\mathbf{P}\|_2$ for (46) poses no practical difficulty. The trade-off between (45) and (46) lies in the relative cost between matrix inversion $\mathbf{A}^{-1} \mathbf{B}$ and multiplication $\mathbf{A} \mathbf{B}$ for square \mathbf{A} and \mathbf{B} . (45) costs one inversion while (46) costs three multiplications. In practice, computing $\mathbf{A}^{-1} \mathbf{B}$ is about 1.5 times the cost of computing $\mathbf{A} \mathbf{B}$ (based on testing in MATLAB and [29]), so overall Eqs. (45) and (46) are equally good alternatives to the SVD/EVD approach for computing the SPD square root or the orthogonal polar factor.

Finally, we emphasize that unlike Eqs. (28)–(31) which aim to solve $\mathbf{X}^2 = \mathbf{P}$, Eqs. (45) and (46) target at solving $\mathbf{X}^T \mathbf{X} = \mathbf{P}$. Consequently, Eqs. (45) and (46) will not yield the major square root of a general \mathbf{P} with eigenvalues not on the negative real axis.

The various new algorithms recommended for solving the minimum-correction second-moment matching problem in Theorem 2.1 are summarized in Table 1.

3.3 The special case of orthonormalization

Now we turn to the special case of orthonormalization, i.e. $\tilde{\mathbf{P}} = \mathbf{I}$, where (27) reduces to $\mathbf{A}_* = \mathbf{P}^{-1/2}$ and $\mathbf{U} = \tilde{U}_* \mathbf{A}_*^{-1}$ becomes the polar decomposition. In this case, [28] provides an a priori bound for $\|\tilde{U}_* - \mathbf{U}\|_F$:

$$\|U^T U - I\|_F / (\|U\|_2 + 1) \leq \|\tilde{U}_* - U\|_F \leq \|U^T U - I\|_F. \tag{51}$$

When $\|U\|_2$ is of order 1 or smaller, the above yields a tight estimate of the order of magnitude of $\|\tilde{U}_* - U\|_F$. One common scenario of such is with U being nearly orthogonal, i.e. $P = U^T U \approx I$, when, for example, a path on a Stiefel manifold is tracked numerically for subspace tracking or optimization with orthogonality constraints.

For orthonormalization, there are many widely-used algorithms that yield an orthogonal matrix not being the closest to the original U . We clarify next how they are related to the minimum-correction solution.

3.3.1 Comparing with QR-based orthonormalization

If we only seek an orthogonal \tilde{U} close to U , but not necessarily the closest one, we can stop after the QR factorization step $U = VR$ in Sect. 3.1 and take $\tilde{U} = V$. This corresponds to the Cholesky decomposition $P = R^T R$. The continuity of the Cholesky factor R with respect to P [47, sec. 12.1.3] guarantees that $P \approx I$ implies $R \approx I$ and thus $V \approx U$. More precisely,

$$\|\tilde{U}_* - U\|_F \leq \|V - U\|_F \leq \frac{1 + \|U\|_2}{\sqrt{2}(1 - \|U^T U - I\|_2)} \|\tilde{U}_* - U\|_F, \tag{52}$$

which is credited to Jiguang Sun by [25], bounds the closeness of the orthogonal QR factor V to U by a multiple of the minimum $\|\tilde{U}_* - U\|_F$ (see also [9, sec. 4]). In particular, we have approximately $\|V - U\|_F \leq \sqrt{2}\|\tilde{U}_* - U\|_F$ when $P \approx I$.

This closeness property of the QR factorization stems from the nested subspace correspondence between V and U , meaning that the first $k = 1, \dots, n$ columns of V and U span the same subspace in \mathbb{R}^m , which constrains V to be close to U . This also implies that when U 's columns are permuted, the corresponding orthogonal QR factor V cannot be obtained by applying the same permutation to the original one, unlike the orthogonal polar factor \tilde{U}_* .

3.3.2 Comparing with plain SVD orthonormalization

[52, app. B] proposes a plain SVD algorithm where the left singular vectors V of an SVD $U = V \Sigma Z^T$ is directly taken as \tilde{U} . The prime issue with this approach is the non-uniqueness of an SVD. The singular vectors associated with the same (or numerically close) singular value are only determined up to an orthogonal transform, which can flip signs, permute the order and even rotate the vectors. This is especially problematic when U is nearly orthogonal because all its singular values will cluster around 1. In this case, $\|V - U\|_F$ can be of order 1 even if U is nearly orthogonal. To remedy this issue, [18] suggest a gradient descent algorithm for orthonormalization. Here we analyze it in ‘‘Appendix D’’ and conclude that due to its linear convergence and stringent stability constraint, it is unattractive compared to the algorithms in Sect. 3.2.

Table 2 Description of the various orthonormalization algorithms benchmarked. Here U is to be orthonormalized and $P = U^T U$ is its Gram matrix

Algorithm	Description
SVD-EVD	Plain SVD ortho. (see Sect. 3.3) by EVD of P
PD-EVD	PD ortho. by EVD of P (see Sect. 3.2.1)
PD-SVD*	PD ortho. by SVD of U (see Sect. 3.2.1)
PD-NtSqr	PD ortho. by computing $P^{1/2}$ using Newton iteration (45)
PD-NtISqr	PD ortho. by computing $P^{-1/2}$ using Newton iteration (46)
PD-GD	PD ortho. by computing $P^{-1/2}$ using gradient descent (87)
PD-NtPD*	PD ortho. by Householder QR $U = VR$ and computing the orthogonal polar factor of R using Newton iteration (45)
QR-CGS*	QR ortho. by classical Gram-Schmidt
QR-MGS*	QR ortho. by modified Gram-Schmidt
QR-Chol	QR ortho. by Cholesky decomposition of P
QR-HH*	QR ortho. by Householder QR $U = VR$

The algorithms that do not form P and directly operate on U are marked by a “*” in their names

4 Benchmarks of the algorithms' performance

In this section, we benchmark the performance of the various new minimum-correction second-moment matching algorithms introduced in Sect. 3 for the special case of orthonormalization (i.e. $P = I$) and compare them to existing algorithms. All the algorithms to be tested are described in Table 2. Note that PD-NtSqr, PD-NtISqr and PD-NtPD are based on our new multi-purpose Newton iteration (45) and (46), while others are existing algorithms from literature. Comparing Table 2 to Table 1, here PD-EVD corresponds to A_* -EVD in Table 1, PD-NtSqr and PD-NtISqr to A_* -NtSqr, while PD-NtPD corresponds to \tilde{U}_* -NtPD.

For all algorithms, the fixed-point iteration terminates once $\|X_k - X_{k-1}\| < 10^{-14} \|X_{k-1}\|$ is satisfied or the maximum number of iterations is reached. The numerical tests are performed under the MATLAB[®] environment on a machine with an Intel[®] Core[™] i7-4702MQ processor and 16GB of memory.

The algorithms are tested in three cases. The first is with a $10^6 \times 10^2$ matrix U (m -by- n , $m \gg n$) of a small condition number $\kappa = 1.5$, so the tall thin U is nearly orthogonal. We form $U = U_0 A Q_0$ by random generation of an orthogonal m -by- n U_0 and orthogonal n -by- n Q_0 . A random orthogonal matrix is obtained by applying QR Householder to a random matrix whose entries are independent and uniformly distributed over $[-1, 1]$. A is a diagonal matrix of condition number κ and of diagonal entries decaying geometrically to 1. Each algorithm is tested using the same 10 samples of U and the average performance metrics are listed in Table 3. For the second and third cases, the same tests are performed, but for a square well-conditioned U and a rectangular ill-conditioned U , respectively. Results are in Table 4.

Overall, Tables 3 and 4 confirm several key messages:

Table 3 Performance of the various algorithms for orthonormalization in the case of well-conditioned (nearly orthogonal) rectangular matrices (m -by- n , $m \gg n$)

	$m = 10^6, n = 10^2, \kappa(U) = 1.5$			
	$\ \tilde{U}^T \tilde{U} - I\ $	$\ \tilde{U} - U\ $	$\ A - A^T\ $	(#Ite) time
SVD-EVD	2e-14	16	12	0.6s
PD-EVD	3e-14	2.75	1e-15	0.6s
PD-SVD*	3e-14	2.75	1e-15	5s
PD-NtSqr	8e-15	2.75	8e-16	(6) 2.3s
PD-NtISqr	9e-15	2.75	1e-15	(8) 0.6s
PD-GD	8e-13	2.75	3e-15	(131) 0.6s
PD-NtPD*	1e-14	2.75	1e-15	(7) 5s
QR-CGS*	6e-15	3.04	1.9	38s
QR-MGS*	6e-15	3.04	1.9	38s
QR-Chol	7e-15	3.04	1.9	2s
QR-HH*	1e-14	3.04	1.9	4s

Here A is the linear transform such that $\tilde{U} = UA$. The Frobenius norm is exclusively used. For the algorithms based on a fixed-point iteration, the average number of iterations is recorded in the parenthesis in front of the run time. In this case, we have $\|U^T U - I\| = 6.5$ for all 10 repetitions, which can be compared to $\|\tilde{U}^T \tilde{U} - I\|$

Table 4 Performance of the algorithms on well-conditioned square matrices ($\kappa(U) = 1.5, \|U^T U - I\| = 29$) and on ill-conditioned rectangular matrices ($\kappa(U) = 10^6, \|U^T U - I\| = 1.5 \times 10^{12}$), respectively

	$m = n = 2 \times 10^3, \kappa(U) = 1.5$			$m = 10^6, n = 10^2, \kappa = 10^6$	
	$\ \tilde{P} - I\ $	$\ \tilde{U} - U\ $	(#Ite) time	$\ \tilde{P} - I\ $	(#Ite) time
PD-EVD	3e-13	12.3	2s	4e-5	0.6s
PD-SVD*	4e-13	12.3	4s	3e-14	5s
PD-NtSqr	4e-14	12.3	(7) 3s	4e-5	(100) 2s
PD-NtInv	6e-14	12.3	(9) 6s	4e-5	(100) 0.6s
PD-GD	3e-11	12.3	(422) 280s	9	(1e4) 3s
PD-NtPD*	7e-14	12.3	(9) 4s	2e-14	(26) 5s
QR-CGS*	8e-14	13.6	13s	6e-5	38s
QR-MGS*	6e-14	13.6	13s	2e-10	38s
QR-Chol	3e-14	13.6	0.3s	3e-5	2s
QR-HH*	7e-14	13.6	0.4s	1e-14	4s

Notations are as those in Table 3

i The plain SVD algorithm does not promote closeness between \tilde{U} and U , as is verified by $\|\tilde{U} - U\|$ of SVD-EVD being much larger than that of the QR-based and of the new minimum-correction (PD-based) algorithms, see Table 3.

ii The new PD-based algorithms achieve the minimum correction as is verified by the agreement between their $\|\tilde{U} - U\|$ values and the analytical one computed using (6) as $\|\tilde{U}_* - U\| = (\text{tr} I + \text{tr} P - 2\text{tr} P^{1/2})^{1/2} = (n + \|U\|_F^2 - 2\|U\|_*)^{1/2}$, where $\|\cdot\|_*$

is the nuclear norm, being the sum of singular values. The symmetry of A_* is verified by the machine epsilon values of $\|A - A^T\|$. The $\|\tilde{U} - U\|$ values of the QR-based algorithms are larger than $\|\tilde{U}_* - U\|$, but within the bound (52).

iii For ill-conditioned matrices, the algorithms based on forming $P = U^T U$ suffer from poor accuracy due to squaring the condition number. This is revealed by the large $\|\tilde{P} - I\|$ values of the algorithms without a “*” in their names, see Table 4. Among them, those based on a fixed-point iteration have their solutions stagnate with large orthogonality errors after only a few iterations. The iterates neither converge nor blow up and just bounce around, even after as many as 10^5 iterations. CGS and MGS are known to suffer from the same orthogonality issue in this case. Hence, the only algorithms that can still reliably orthonormalize U are those based on performing Householder QR or SVD on U .

iv For tall thin matrices, the algorithms based on forming $P = U^T U$ tend to be more efficient when U is well-conditioned, because other than computing the small P in the beginning and applying A to U in the end, they only deal with small n -by- n matrices. This strategy avoids performing QR factorization or SVD to a large m -by- n matrix, which could be relatively expensive in this case.

There is yet another algorithm for computing the SPD matrix square root that takes a completely different route from fixed-point iterations. [23] utilizes the matrix version of the Cauchy integral formula for the square root function and combines conformal maps and the trapezoidal rule in computing the contour integral. However, its exponential/geometric convergence rate is only equivalent to a linearly converging fixed-point iteration, which is significantly slower than the quadratic convergence of Eqs. (45) and (46) proposed previously. Indeed, this can be verified by comparing to the results based on the Pascal matrices shown in Fig. 7 of [23]. It takes 9 and 20 Newton iterations to obtain the SPD square root of the 3-by-3 and 8-by-8 Pascal matrix, respectively, while the contour-integral approach requires 13 and 34 quadrature points to reach the same accuracy. Here every extra quadrature point has the equivalent cost of one iteration step because both involve solving a linear system of the same size. Besides efficiency, using fixed-point iteration has the extra advantage of conceptual and algorithmic simplicity.

5 Application to matrix differential equations that preserve orthogonality

A matrix differential equation

$$d_t U = F(U) \quad (53)$$

for $U(t) \in \mathcal{M}_{m \times n}$ is said to preserve orthogonality if $U^T U = I$ at $t = 0$ implies $U^T U = I$ at all $t > 0$. We encounter such equations in applications such as subspace tracking [5,8,17,30,32,34], where we track the time evolution of an n -dimensional subspace of \mathbb{R}^m , which may come from a discretization of an infinite-dimensional dynamical system, by evolving an orthonormal basis with its base vectors forming U 's columns. In this section, we showcase our new minimum-correction second-moment

matching algorithms in the special case of $\tilde{P} = I$ for re-orthonormalization required in solving such equations. Before the examples, we analyze the error due to the loss of numerical orthogonality and clarify the role played by re-orthonormalization.

Although (53) preserves orthogonality, its time-discretized counterpart $\hat{U}_{n+1} = \hat{F}(\hat{U}_n)$ in general does not [11,24,25,53]. Here, \hat{U}_n is a numerical approximation of $U_n = U(t_n)$ and the map \hat{F} stems from a particular time-marching scheme, which is said to preserve orthogonality if $\hat{U}_n^T \hat{U}_n = I \implies \hat{U}_{n+1}^T \hat{U}_{n+1} = I$. According to [11,24], the only known family of orthogonality-preserving time-marching schemes is the one of Gauss–Legendre Runge–Kutta schemes. Unfortunately, such schemes are difficult to implement, expensive to use and not widely available in numerical software [25]. Moreover, they only preserve orthogonality for certain F , which renders them even less attractive. Therefore, a more practical strategy is to integrate (53) by an arbitrary scheme and re-orthonormalize the \hat{U}_n 's at the end of each time step.

We denote by \hat{U} the numerical solution at a particular time obtained by some time-marching scheme without re-orthonormalization, and by $\Delta U = \hat{U} - U$ the time-marching error. If we choose $\|\cdot\|$ as the matrix 2-norm, we have

$$\|\hat{U}^T \hat{U} - I\| = \|(U + \Delta U)^T (U + \Delta U) - I\| \leq 2 \|\Delta U\| \|U\| + \|\Delta U\|^2,$$

since $U^T U = I$. If the scheme is p -th order accurate, i.e. $\|\Delta U\| = O((\Delta t)^p)$, since $\|U\| = O(1)$, we have $\|\hat{U}^T \hat{U} - I\| = O((\Delta t)^p)$, which is of the same order as the time-marching error. This seems to imply that as long as the time-marching scheme is accurate enough, there is no need to worry about the orthogonality error. However, the fact that (53) preserves orthogonality only implies that the Stiefel manifold $S_I = \{U : U^T U = I\}$ is invariant under (53), but S_I may not be a stable manifold. Therefore, any deviation of U from S_I may not stay bounded under (53). It might even hit an unstable direction and diverge exponentially. In addition, in many cases, the derivation and thus the validity of (53) itself relies on U remaining orthogonal, so a violation of this condition may unexpectedly ruin the numerical solution.

The above analysis justifies re-orthonormalizing \hat{U} after each time step. This procedure should ideally not undermine the convergence order of the time-marching scheme or introduce significant numerical artifacts. Hence, we need to analyze how re-orthonormalization affects the numerical errors.

Suppose $\tilde{U} \in S_I$ is the orthonormalized solution obtained from \hat{U} . To maintain the order of convergence, it suffices to have $\|\tilde{U} - \hat{U}\| \leq C \|\hat{U} - U\|$ for some $C > 0$ because then by triangle inequality, we have $\|\tilde{U} - U\| \leq (C + 1) \|\hat{U} - U\|$. This can be satisfied if we choose \tilde{U} to be \tilde{U}_* , the orthogonal projection of \hat{U} onto S_I , because, as $U \in S_I$, Corollary 2.3 implies $\|\tilde{U}_* - \hat{U}\| \leq \|\hat{U} - U\|$ and thus $\|\tilde{U}_* - U\| \leq 2 \|\hat{U} - U\|$. Indeed, since $\|\hat{U} - U\|$ is small, \tilde{U}_* , \hat{U} , and U are almost the vertices of a right triangle (“almost” because S_I is not flat), so in practice, $\|\tilde{U}_* - U\| \leq \|\hat{U} - U\|$. Numerical examples in [11] show that in some cases $\|\tilde{U}_* - U\|$ can even be smaller than $\|\hat{U} - U\|$ by an order of magnitude. Another viable option for \tilde{U} is the orthogonal QR factor V of \hat{U} because (52) guarantees $\|V - \hat{U}\| \leq \sqrt{2} \|\tilde{U}_* - \hat{U}\| \leq \sqrt{2} \|\hat{U} - U\|$.

To showcase the new algorithms and demonstrate the above results, we consider the stochastic Lorenz-96 system [38] (the index i is circular, i.e. $x_i = x_{i+m}$),

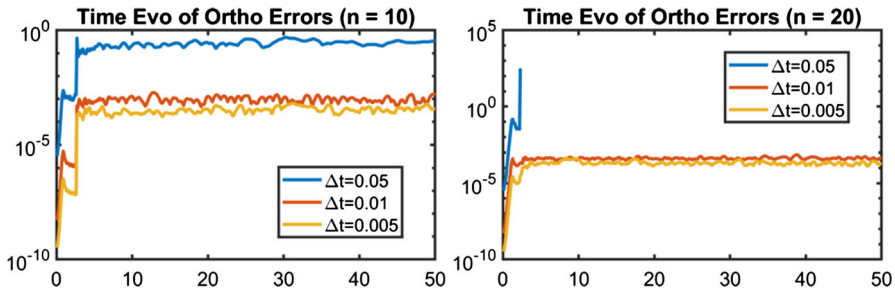


Fig. 2 Orthogonality errors in the absence of re-orthonormalization. The error saturates at a higher level with larger time step Δt . When $\Delta t = 0.05$, the orthogonality is completely lost for $n = 10$ (left), while the solution even blows up for $n = 20$ (right)

$$d_t x_i = (x_{i+1} - x_{i-2})x_{i-1} - x_i + F, \quad i = 1, \dots, m. \tag{54}$$

This toy model mimics atmospheric convection along a mid-latitude circle. We set $m = 40$ and $F = 8$ which lead to chaotic dynamics. In all simulations, the initial condition is $x_i = 0$ for $i \neq 1$ while x_1 is uniformly distributed over $[-0.01, 0.01]$. The classical explicit RK-4 scheme is used for time integration.

To solve (54), we employ the dynamically orthogonal (DO) equations [46]. They approximate the random vector $\mathbf{u}(t; \omega) = [x_1, \dots, x_m]^T \in \mathbb{R}^m$ by a low rank expansion $\mathbf{u}(t; \omega) = \bar{\mathbf{u}}(t) + \sum_{i=1}^n \phi_i(t; \omega) \mathbf{u}_i(t) = \bar{\mathbf{u}} + \mathbf{U}\boldsymbol{\phi}$ where $\bar{\mathbf{u}}$ is the mean, $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_n] \in \mathcal{M}_{m \times n}$ contains the orthonormal base vectors and $\boldsymbol{\phi} = [\phi_1, \dots, \phi_n]^T$ contains the random scalar coefficients. The evolution of all these components is governed by a dynamical system, the DO equations, that preserves and whose validity is predicated on $\mathbf{U}^T \mathbf{U} \equiv \mathbf{I}$. Such techniques are useful for uncertainty quantification in high-dimensional systems governed by time-dependent stochastic PDEs [18,49]. Related works include the dynamic low-rank approximation of a time-varying matrix [30], geometric analyses of the DO equations based on matrix manifold theory [17,19] and dynamically bi-orthogonal equations [10].

In Fig. 2, we first show the time evolution of the orthogonality error $\|\mathbf{U}^T \mathbf{U} - \mathbf{I}\|_{\max}$ ($\|A\|_{\max} \triangleq \max\{A_{ij}\}$) for the numerical integration of the DO equations without any re-orthonormalization. Results are shown for $n = 10, 20$ and three time step sizes. As we can see, the orthogonality error grows with Δt since a larger Δt implies a larger discrepancy between the time-discrete system and the originally orthogonality-preserving continuous one. Moreover, when $\Delta t = 0.005$ and 0.01 , $\|\mathbf{U}^T \mathbf{U} - \mathbf{I}\|_{\max}$ tends to saturate at a low level ($< 10^{-3}$) that corresponds to only a negligible deviation from orthogonality. However, when Δt increases to 0.05 , $\|\mathbf{U}^T \mathbf{U} - \mathbf{I}\|_{\max}$ can either level off at a magnitude of 1 for $n = 10$ that signals a complete loss of orthogonality or even drive the simulation unstable and cause the solution to blow up for $n = 20$. This exemplifies the necessity of re-orthonormalization.

Note that in the expansion $\mathbf{u} = \bar{\mathbf{u}} + \mathbf{U}\boldsymbol{\phi}$, there is also the unitary freedom $\mathbf{U}\boldsymbol{\phi} = (\mathbf{U}\mathbf{Q})(\mathbf{Q}^T\boldsymbol{\phi})$ for any orthogonal \mathbf{Q} , under which all realizations $\mathbf{u}(t; \omega)$ are unchanged. Therefore, if we apply such a \mathbf{Q} at some time t and use $\mathbf{U}\mathbf{Q}$ and $\mathbf{Q}^T\boldsymbol{\phi}$ as the initial values of the modes and coefficients for the subsequent time integra-

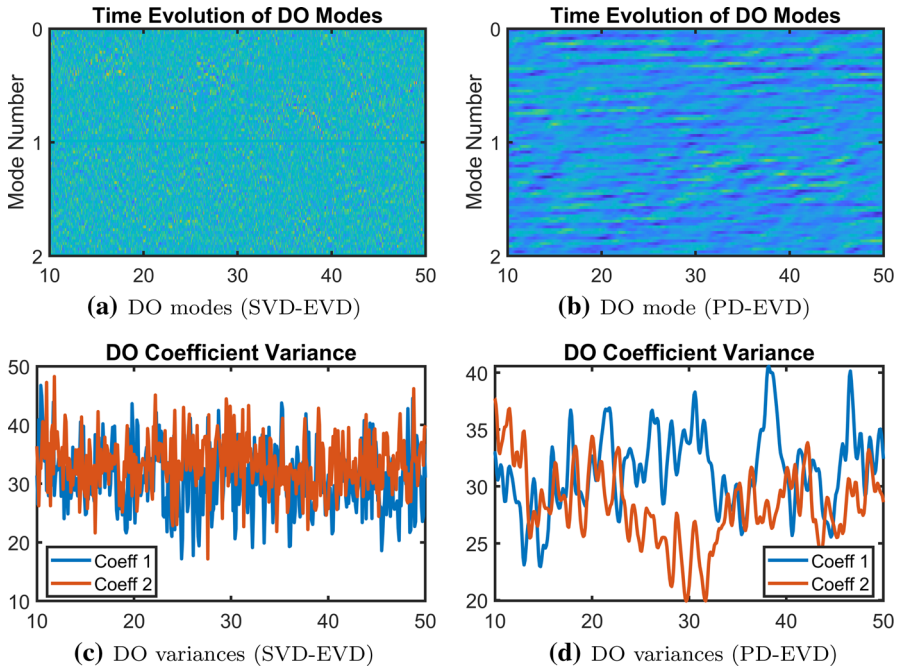


Fig. 3 Comparison between orthonormalization using SVD–EVD (left) with that using polar decomposition PD-EVD (right). The evolution (from $t = 10$ to 50 , $\Delta t = 0.05$) of the first two DO modes $[\mathbf{u}_1^T, \mathbf{u}_2^T]^T$ stacked together are shown on the top and the evolution of their coefficient variances $\text{Var}(\phi_1)$ and $\text{Var}(\phi_2)$ on the bottom

tion, the future evolution of $\mathbf{u}(t; \omega)$ will remain the same. However, this breaks the time continuity of \mathbf{U} and ϕ at time t , as illustrated in Fig. 3, where the SVD–EVD and PD-EVD (see Table 2) algorithm are compared. Due to the non-uniqueness of an SVD, SVD–EVD not only removes the orthogonality error, but also applies to \mathbf{U} some \mathbf{Q} that may not at all be close to \mathbf{I} . As a result, we lose the time continuity of the modes and coefficients as indicated by the “random” pattern in the left half of Fig. 3. In contrast, PD-EVD minimizes such numerical artifacts. Maintaining time continuity is especially important when we use multi-step time-marching schemes⁵ because their validity relies on the smooth evolution of \mathbf{U} and ϕ across different time steps. A “random” \mathbf{Q} applied at each time step destroys this smoothness and can cause the time integration to diverge.

Finally we show in Fig. 4 how the orthogonality error can be controlled by PD-EVD, PD-NtSqr, PD-NtISqr, QR-Chol and SVD–EVD for re-orthonormalization. These algorithms were described in Table 2. They are the best choices in this case based on the discussion in Sect. 4 because the $\hat{\mathbf{U}}$ to be re-orthonormalized after each time step is a tall thin matrix and is already nearly orthogonal. Here we exclusively use $\Delta t = 0.05$ and $n = 10$ for all simulations. The left of Fig. 4 indicates that the orthogonality errors are kept at almost machine epsilon. The two EVD-based algorithms have slightly larger

⁵ For example, [49] uses the leapfrog scheme and [13] uses the implicit-explicit backward difference schemes for the DO equations.

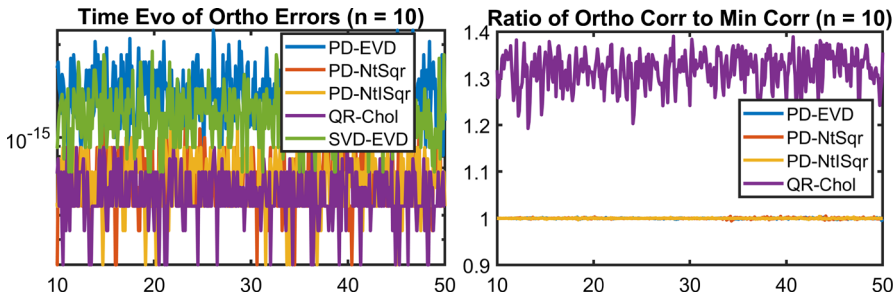


Fig. 4 The orthogonality errors after re-orthonormalization by various algorithms (left) and the ratio of the correction $\|\tilde{\mathbf{U}} - \hat{\mathbf{U}}\|$ to the minimum $\|\tilde{\mathbf{U}}_* - \hat{\mathbf{U}}\|$ computed by (6) (right). The ratios for SVD–EVD are of order 10^6 and thus not plotted. Here $\Delta t = 0.05$ and $n = 10$ are used for all the simulations

errors, most likely due to the convergence criterion for the built-in EVD subroutine, but this is not essential in practice. The right plot shows the ratio of the actual correction $\|\tilde{\mathbf{U}} - \hat{\mathbf{U}}\|$ to the theoretical minimum $\|\tilde{\mathbf{U}}_* - \hat{\mathbf{U}}\| = (n + \text{tr} \mathbf{P} - 2\text{tr} \mathbf{P}^{1/2})^{1/2}$ computed by (6). As we can see, all three PD-based algorithms achieve a ratio of 1, while the QR-based algorithm has a ratio between 1.2 and 1.4, agreeing with the bound (52). Not surprisingly, the plain SVD orthonormalization yields a huge ratio of order 10^6 (not plotted) because it produces corrections of order 1 even when $\|\hat{\mathbf{U}}^T \tilde{\mathbf{U}} - \mathbf{I}\|$ is only of order 10^{-6} .

6 Application to ensemble square root filters for data assimilation

Another application of the minimum-correction second-moment matching is in the ensemble square root filters for data assimilation [3,6,15,16,35,41,44]. Given an n -dimensional dynamical system $d_t \mathbf{u} = \mathbf{f}(\mathbf{u})$ as well as a linear observation model $\mathbf{y} = \mathbf{H}\mathbf{u} + \mathbf{e} \in \mathbb{R}^d$ with Gaussian noise $\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$, the task of data assimilation is to combine the information content of observations $\hat{\mathbf{y}}$ with that of model predictions to refine the (probabilistic) estimate of the state \mathbf{u} . At each observation time t , we denote a forecast ensemble with m realizations for \mathbf{u} by $\mathbf{U}_- \in \mathcal{M}_{n \times m}$ ⁶ and its sample covariance by $\mathbf{P}_- = (1/m)\tilde{\mathbf{U}}_- \tilde{\mathbf{U}}_-^T$. Here, \mathbf{U}_- is decomposed into the mean and the fluctuation as $\mathbf{U}_- = \bar{\mathbf{u}}_- \mathbf{1}^T + \tilde{\mathbf{U}}_-$, where $\mathbf{1}$ is a vector with all entries being 1.

The task of filtering, the most common type of data assimilation, is to obtain an analysis ensemble $\mathbf{U}_+ \in \mathcal{M}_{n \times m}$, which integrates information in \mathbf{U}_- with the observations collected at t denoted by $\hat{\mathbf{y}}$, so that \mathbf{U}_+ better captures the underlying true state than \mathbf{U}_- can. A popular family of filters is that of ensemble square root filters, which update the mean in the same ways as a Kalman filter:

$$\bar{\mathbf{u}}_+ = \bar{\mathbf{u}}_- + \mathbf{K}(\hat{\mathbf{y}} - \mathbf{H}\bar{\mathbf{u}}_-), \quad \mathbf{K} = \mathbf{P}_- \mathbf{H}^T (\mathbf{H} \mathbf{P}_- \mathbf{H}^T + \mathbf{R})^{-1}, \quad (55)$$

⁶ Here the role of rows and columns flips compared to previous sections, since we want to follow the notation convention in the field of data assimilation.

and update the ensemble spread in some way such that the sample covariance of U_+ matches the one given by a Kalman filter, i.e.

$$P_+ = (1/m)\tilde{U}_+\tilde{U}_+^T = (I - KH)P_- \tag{56}$$

Many variants for obtaining such a \tilde{U}_+ have been proposed, including the ensemble Kalman filter [15,16], error subspace statistical estimation [33,35], ensemble adjustment filter [3] and ensemble transform filter [6], to name a few [54]. Here, we demonstrate a particular variant proposed by [41], which choose \tilde{U}_+ to be the one closest to \tilde{U}_- , reducing the task to minimum-correction second-moment matching. This choice may be seen as retaining as much physical information from the prior ensemble U_- as possible [41].

More precisely, the update formula for \tilde{U}_+ proposed by [41] is

$$\tilde{U}_+ = \underset{V: VV^T = mP_+}{\operatorname{arg\,min}} \|V - \tilde{U}_-\|_{F, P_-^{-1}} \tag{57}$$

Corollary 2.1 implies that $\tilde{U}_+ = A\tilde{U}_-$ with (here $\Gamma_n = P_-^{-1}$)

$$A = \sqrt{\Gamma_n}^{-1} A_* \sqrt{\Gamma_n} = \sqrt{P_-^{-1}} (\sqrt{P_-^{-1}} P_+ \sqrt{P_-^{-1}})^{1/2} \sqrt{P_-^{-1}} \tag{58}$$

This is almost the same as $A = P_-^{1/2} (P_-^{-1/2} P_+ P_-^{-1/2})^{1/2} P_-^{-1/2}$, which is eq. A14 in [41], except that all the square root $\sqrt{\cdot}$ were unnecessarily restricted to the unique SPD one $(\cdot)^{1/2}$ in the expression used by [41].

We will again use the Lorenz-96 system (54) as an example with exactly the same numerical setups as those in Sect. 5, except that $F = 4$ is used now. Hence, we have $n = 40, m = 1000$. Note that n now denotes the system’s dimension while m the number of realizations, differing from Sect. 5. The observation data include every other variable, i.e. x_1, x_3, \dots, x_{39} , at time $t = 10, 15, 20, \dots, 50$. The observation noise covariance $R = 0.01I$ is assumed known.

Here four filters are tested and compared: “EnSQR” is the one proposed by [41] and based on (57); “EnKF” is the classical ensemble Kalman filter [16]; “KF” is the same as EnKF except that the filtering update is replaced by fitting a Gaussian to the prior ensemble, applying the classical Kalman update and re-sampling from the posterior Gaussian. These three filters are based on brute force Monte-Carlo simulations of the Lorenz-96 system without any dimension reduction. “EnSQR-DO” is EnSQR applied to a DO simulation with a 10-dimensional reduced subspace. The performance of these four filters are compared in Fig. 5.

The RMSE (root mean squared errors) on the left shows that all filters perform similarly except that EnSQR-DO has a larger error due to dimension reduction (the attractor is not well contained in a 10-dimensional linear subspace globally). The ratios of the ensemble correction (fluctuation part) $\|\tilde{U}_+ - \tilde{U}_-\|$ at each filtering step to the theoretical minimum $\|\tilde{U}_{+*} - \tilde{U}_-\|$ for posterior covariance matching are plotted on the right. The minimum is computed by (12) being $\|\tilde{U}_{+*} - \tilde{U}_-\|^2 = \operatorname{tr}(I + \tilde{P}_\Gamma - 2\tilde{P}_\Gamma^{1/2})$, where $\tilde{P}_\Gamma = \sqrt{P_-^{-1}} P_+ \sqrt{P_-^{-1}}$. As we can see, the two minimum-correction-based

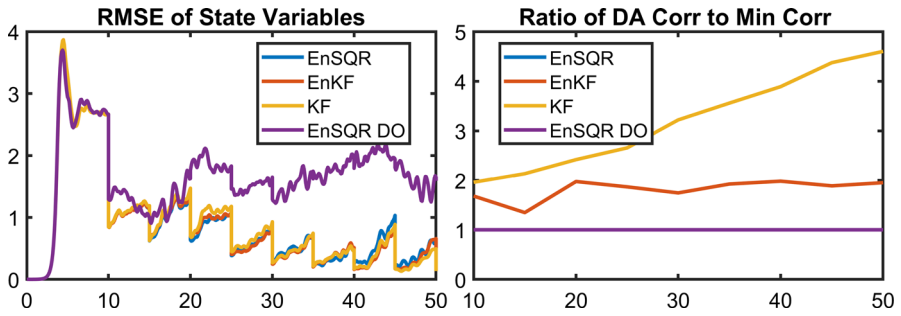


Fig. 5 Time evolution of RMSE of the various filters (left) and the ratio of the filtering correction $\|\tilde{U}_+ - \tilde{U}_-\|$ to the minimum $\|\tilde{U}_{+*} - \tilde{U}_-\|$ computed by (12) (right). The first three cases are without dimension reduction while the last one is reduced to a 10-dimensional subspace using DO. Here $\Delta t = 0.05$ is used for all the simulations

filters (EnSQR, EnSQR-DO) achieve the unit ratio while the other two do not. KF has the largest ratio because its re-sampling step completely breaks any realization correspondence between the prior and posterior ensemble. In contrast, EnKF adopts a realization-wise update strategy, which yields a ratio between those of KF and of the minimum-correction filters. For further details on the performance of the minimum-correction-based ensemble square root filter (EnSQR), please refer to section 5.3.3 and 5.7 of [36].

7 Conclusions

In this paper, we develop the theory for minimum-correction second-moment matching. We solve the optimization analytically for both the full-rank and rank-deficient case. The finite-dimensional result is generalized to the infinite-dimensional setting and a connection to optimal transport is drawn. We also show how the general problem can reduce to the familiar one of optimal orthonormalization and we accordingly generalize the polar decomposition, which is known to solve this special case.

We obtain numerical schemes for computing the optimizer. We show the instrumental role played by the algorithms for polar decomposition and SPD matrix square root and we analyze existing ones in the literature. We modify two Newton iterations deemed unstable before and significantly improve their stability. The resulting two new multi-purpose schemes (45) and (46) can be used to efficiently compute both the orthogonal polar factor and the SPD square root. These iterations play a key role in minimum-correction second-moment matching using algorithm A_{*} -NtSqr and \tilde{U}_{*} -NtPD in Table 1 and in their counterparts PD-NtSqr, PD-NtISqr, and PD-NtPD in Table 2 for orthonormalization. We verify the higher performance of the new algorithms using benchmarks with random matrices.

Finally, we showcase results in two applications. In reduced subspace tracking for uncertainty quantification, we maintain the numerical orthogonality and continuity of a time-varying orthonormal basis. In an ensemble square root filter for data assimilation,

we obtain the unique posterior ensemble that matches a given covariance matrix while also minimizing its distance to the prior one.

Acknowledgements We thank the members of our MSEAS group at MIT. We especially thank Florian Feppon for several suggestions and for pointing us to the connection to the polar decomposition. We are grateful to the Office of Naval Research for support under grants N00014-14-1-0476 (Science of Autonomy – LEARNS) and N00014-14-1-0725 (Bays-DA) to the Massachusetts Institute of Technology. We also want to thank the anonymous referees for reading our manuscript carefully and providing helpful feedback to help us improve the presentation.

A: Proof of theorem 2.1

Proof First, since the objective function of (4) can be rewritten as

$$\|\tilde{U} - U\|_F^2 = \text{tr}[(\tilde{U} - U)^T(\tilde{U} - U)] = \text{tr}\tilde{P} + \text{tr}P - 2\text{tr}(U^T\tilde{U}),$$

we have

$$\arg \min_{\tilde{U} \in \mathcal{M}_{m \times n}: \tilde{U}^T \tilde{U} = \tilde{P}} \|\tilde{U} - U\|_F^2 = \arg \min_{\tilde{U} \in \mathcal{M}_{m \times n}: \tilde{U}^T \tilde{U} = \tilde{P}} -2\text{tr}(U^T\tilde{U}). \tag{59}$$

The Lagrangian of this optimization is thus

$$L(\tilde{U}, \Lambda) = -2\text{tr}(U^T\tilde{U}) + \text{tr}(\Lambda^T(\tilde{U}^T\tilde{U} - \tilde{P})), \tag{60}$$

where $\Lambda \in \mathcal{M}_{n \times n}$ is the Lagrangian multiplier for the constraint $\tilde{U}^T\tilde{U} = \tilde{P}$. The global optimizer should be one of the critical points, so ⁷

$$\nabla_{\tilde{U}} L|_{\tilde{U}_*, \Lambda_*} = -2U + 2\tilde{U}_* \Lambda_* = \mathbf{0}, \quad \nabla_{\Lambda} L|_{\tilde{U}_*, \Lambda_*} = \tilde{U}_*^T \tilde{U}_* - \tilde{P} = \mathbf{0}. \tag{61}$$

The first condition gives $U = \tilde{U}_* \Lambda_*$ and the second is simply the second-moment constraint on \tilde{U}_* . The symmetry of this constraint implies the symmetry of the Lagrangian multiplier Λ_* . Inserting these results into the objective function (59), we obtain

$$-2\text{tr}(U^T\tilde{U}_*) = -2\text{tr}(\Lambda_*^T \tilde{U}_*^T \tilde{U}_*) = -2\text{tr}(\Lambda_* \tilde{P}) = -2\text{tr}(\sqrt{\tilde{P}} \Lambda_* \sqrt{\tilde{P}}^T). \tag{62}$$

The reason for symmetrifying the matrix in the last step will soon become clear. Since

$$P = U^T U = \Lambda_* \tilde{U}_*^T \tilde{U}_* \Lambda_* = \Lambda_* \tilde{P} \Lambda_*,$$

the second-moment constraint reduces to $P = \Lambda_* \tilde{P} \Lambda_*$. To relate the objective function (62) to this constraint, we use a small trick:

⁷ See [42] for how these matrix gradients are computed or refer to ‘‘Appendix A’’ of [37].

$$\begin{aligned}
 (\sqrt{\tilde{\mathbf{P}}}\mathbf{A}_*\sqrt{\tilde{\mathbf{P}}^T})^2 &= (\sqrt{\tilde{\mathbf{P}}}\mathbf{A}_*\sqrt{\tilde{\mathbf{P}}^T})(\sqrt{\tilde{\mathbf{P}}}\mathbf{A}_*\sqrt{\tilde{\mathbf{P}}^T}) \\
 &= \sqrt{\tilde{\mathbf{P}}}(\mathbf{A}_*\tilde{\mathbf{P}}\mathbf{A}_*)\sqrt{\tilde{\mathbf{P}}^T} \\
 &= \sqrt{\tilde{\mathbf{P}}}\mathbf{P}\sqrt{\tilde{\mathbf{P}}^T}.
 \end{aligned}$$

Therefore, $\sqrt{\tilde{\mathbf{P}}}\mathbf{A}_*\sqrt{\tilde{\mathbf{P}}^T}$ must be a symmetric square root of $\sqrt{\tilde{\mathbf{P}}}\mathbf{P}\sqrt{\tilde{\mathbf{P}}^T}$, which is not unique. This characterizes all the critical points of (4), among which is the desired global optimizer $\tilde{\mathbf{U}}_* = \mathbf{U}\mathbf{A}_*^{-1}$.

Next, we identify the global optimizer by comparing the values of the objective function evaluated at these critical points. Since $\sqrt{\tilde{\mathbf{P}}}\mathbf{A}_*\sqrt{\tilde{\mathbf{P}}^T}$ is symmetric, it has an eigendecomposition

$$\sqrt{\tilde{\mathbf{P}}}\mathbf{A}_*\sqrt{\tilde{\mathbf{P}}^T} = \mathbf{V}\text{diag}(\lambda_1, \dots, \lambda_n)\mathbf{V}^T$$

and its square $\sqrt{\tilde{\mathbf{P}}}\mathbf{P}\sqrt{\tilde{\mathbf{P}}^T}$ will have the corresponding eigendecomposition

$$\sqrt{\tilde{\mathbf{P}}}\mathbf{P}\sqrt{\tilde{\mathbf{P}}^T} = \mathbf{V}\text{diag}(\lambda_1^2, \dots, \lambda_n^2)\mathbf{V}^T. \tag{63}$$

If we denote the n positive eigenvalues of $\sqrt{\tilde{\mathbf{P}}}\mathbf{P}\sqrt{\tilde{\mathbf{P}}^T}$ by $\sigma_1, \dots, \sigma_n$, then without loss of generality, we must have $\sigma_i = \lambda_i^2$ and hence $\lambda_i = \pm\sqrt{\sigma_i}$. Therefore, the objective function (62) reduces to

$$\begin{aligned}
 -2\text{tr}(\mathbf{U}^T\tilde{\mathbf{U}}_*) &= -2\text{tr}(\mathbf{V}\text{diag}(\lambda_1, \dots, \lambda_n)\mathbf{V}^T) \\
 &= -2\sum_{i=1}^n \pm\sqrt{\sigma_i} \geq -2\sum_{i=1}^n \sqrt{\sigma_i}.
 \end{aligned} \tag{64}$$

We can see that the lower bound is attained if and only if $\lambda_i = \sqrt{\sigma_i}, i = 1, \dots, n$, i.e. when $\sqrt{\tilde{\mathbf{P}}}\mathbf{A}_*\sqrt{\tilde{\mathbf{P}}^T}$ takes the unique SPD square root:

$$\sqrt{\tilde{\mathbf{P}}}\mathbf{A}_*\sqrt{\tilde{\mathbf{P}}^T} = (\sqrt{\tilde{\mathbf{P}}}\mathbf{P}\sqrt{\tilde{\mathbf{P}}^T})^{1/2}. \tag{65}$$

Therefore, the global minimizer to (4) is $\tilde{\mathbf{U}}_* = \mathbf{U}\mathbf{A}_*^{-1}$ with

$$\mathbf{A}_*^{-1} = \sqrt{\tilde{\mathbf{P}}^T}(\sqrt{\tilde{\mathbf{P}}}\mathbf{P}\sqrt{\tilde{\mathbf{P}}^T})^{-1/2}\sqrt{\tilde{\mathbf{P}}}. \tag{66}$$

Note that the particular choice of square root $\sqrt{\tilde{\mathbf{P}}}$ is irrelevant because both the spectrum of $\sqrt{\tilde{\mathbf{P}}}\mathbf{P}\sqrt{\tilde{\mathbf{P}}^T}$ and the global optimizer (66) are invariant under the unitary freedom $\sqrt{\tilde{\mathbf{P}}} \rightarrow \mathbf{Q}\sqrt{\tilde{\mathbf{P}}}$, due to the fact that $\mathbf{Q}\mathbf{A}^{1/2}\mathbf{Q}^T = (\mathbf{Q}\mathbf{A}\mathbf{Q}^T)^{1/2}$ for any SPD \mathbf{A} .

We remark that an alternative to using (62) is to use

$$-2\text{tr}(\mathbf{U}^T \tilde{\mathbf{U}}_*) = -2\text{tr}(\mathbf{U}^T \mathbf{U} \mathbf{A}_*^{-1}) = -2\text{tr}(\mathbf{P} \mathbf{A}_*^{-1}) = -2\text{tr}(\sqrt{\mathbf{P}} \mathbf{A}_*^{-1} \sqrt{\mathbf{P}}^T)$$

and correspondingly $\tilde{\mathbf{P}} = \mathbf{A}_*^{-1} \mathbf{P} \mathbf{A}_*^{-1}$. This simply switches the role of \mathbf{P} and $\tilde{\mathbf{P}}$ and replaces \mathbf{A}_* with \mathbf{A}_*^{-1} . Therefore, this gives an alternative expression of \mathbf{A}_*^{-1} obtained in (66), i.e.

$$\mathbf{A}_*^{-1} = \sqrt{\mathbf{P}}^{-1} (\sqrt{\mathbf{P}} \tilde{\mathbf{P}} \sqrt{\mathbf{P}}^T)^{1/2} \sqrt{\mathbf{P}}^{-T}. \tag{67}$$

□

B: Generalization from \mathbb{R}^m to a Hilbert space

In Theorem 2.1, each column of \mathbf{U} or $\tilde{\mathbf{U}}$ is an element in \mathbb{R}^m and the second-moment matrices are computed based on an inner product on \mathbb{R}^m . A natural question is whether Theorem 2.1 still holds when \mathbb{R}^m is replaced by an infinite-dimensional Hilbert space \mathcal{H} equipped with an inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$.

More precisely, assume $u_1, \dots, u_n \in \mathcal{H}$ are linearly independent and form an n -vector of \mathcal{H} denoted by $\mathbf{u} = [u_1, \dots, u_n]^T \in \mathcal{H}^n$. Then we can compute the second-moment matrix (a.k.a. the Gram matrix) of \mathbf{u} as $\mathbf{P} = \text{Gr}(\mathbf{u}, \mathbf{u})$, where

$$\text{Gr}(\mathbf{u}, \mathbf{v}) \triangleq \left\langle \mathbf{u}, \mathbf{v}^T \right\rangle_{\mathcal{H}} \in \mathcal{M}_{n \times n} \tag{68}$$

i.e. with the (i, j) -th element of $\text{Gr}(\mathbf{u}, \mathbf{v})$ being $\langle u_i, v_j \rangle_{\mathcal{H}}$. If \mathcal{H}^n is equipped with an inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}^n}$ defined as

$$\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{H}^n} \triangleq \text{tr} \left(\sqrt{\mathbf{\Gamma}_n} \text{Gr}(\mathbf{u}, \mathbf{v}) \sqrt{\mathbf{\Gamma}_n}^T \right) \tag{69}$$

with again an SPD weight matrix $\mathbf{\Gamma}_n$, we want to solve the following minimum-correction second-moment matching problem

$$\arg \min_{\tilde{\mathbf{u}} \in \mathcal{H}^n : \text{Gr}(\tilde{\mathbf{u}}, \tilde{\mathbf{u}}) = \tilde{\mathbf{P}}} \|\tilde{\mathbf{u}} - \mathbf{u}\|_{\mathcal{H}^n}^2 \tag{70}$$

for a given SPD $\tilde{\mathbf{P}} \in \mathcal{M}_{n \times n}$. Here the norm is induced by the inner product (69). As a common scenario, if $\mathcal{H} = L_2(\Omega)$, which is the space of square-integrable real functions on some measure space (Ω, \mathcal{A}, m) equipped with inner product $\langle u(\cdot), v(\cdot) \rangle_{\mathcal{H}} = \int_{\Omega} u(\omega)v(\omega)dm(\omega)$, then the inner product (69) on \mathcal{H}^n can be rewritten as

$$\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{H}^n} = \int_{\Omega} \mathbf{u}(\omega)^T \mathbf{\Gamma}_n \mathbf{v}(\omega) dm(\omega). \tag{71}$$

The definitions (69) and (71) are simply the Hilbert space counterparts of the two expressions in (2). Since we can get rid of $\mathbf{\Gamma}_n$ by scaling the elements in \mathcal{H}^n by $\sqrt{\mathbf{\Gamma}_n}$ as before, we will assume $\mathbf{\Gamma}_n = \mathbf{I}$ hereafter.

Theorem 2.1 does not apply to (70), but as in Sect. 2.2, if we restrict the candidate set to $\{\tilde{\mathbf{u}} = \mathbf{A}\mathbf{u} : \mathbf{A} \in \mathcal{M}_{n \times n}\}$, the Hilbert space complication will be confined to computing the Gram matrices only. Since

$$\text{Gr}(\mathbf{A}\mathbf{u}, \mathbf{B}\mathbf{v}) = \left\langle \mathbf{A}\mathbf{u}, (\mathbf{B}\mathbf{v})^\top \right\rangle_{\mathcal{H}} = \mathbf{A} \left\langle \mathbf{u}, \mathbf{v}^\top \right\rangle_{\mathcal{H}} \mathbf{B}^\top = \mathbf{A}\text{Gr}(\mathbf{u}, \mathbf{v})\mathbf{B}^\top$$

due to the bilinearity of an inner product, we have

$$\|\tilde{\mathbf{u}} - \mathbf{u}\|_{\mathcal{H}^n}^2 = \|(\mathbf{A} - \mathbf{I})\mathbf{u}\|_{\mathcal{H}^n}^2 = \text{tr} \left((\mathbf{A} - \mathbf{I})\text{Gr}(\mathbf{u}, \mathbf{u})(\mathbf{A} - \mathbf{I})^\top \right) = \|\mathbf{A} - \mathbf{I}\|_{\mathbb{F}, \mathbf{P}}^2.$$

Hence (70) reduces to (16) and Corollary 2.2 applies.

To establish the analog of Theorem 2.1 for (70), we need some techniques from the calculus of variation to generalize optimization in \mathbb{R}^m to that in a Hilbert space.

Theorem B.1 (Minimum-correction 2nd-moment matching on a Hilbert space) *Given $\mathbf{u} \in \mathcal{H}^n$ whose entries are linearly independent and $\tilde{\mathbf{P}}$ which is SPD, we have*

$$\arg \min_{\tilde{\mathbf{u}} \in \mathcal{H}^n : \text{Gr}(\tilde{\mathbf{u}}, \tilde{\mathbf{u}}) = \tilde{\mathbf{P}}} \|\tilde{\mathbf{u}} - \mathbf{u}\|_{\mathcal{H}^n}^2 = \mathbf{A}_* \mathbf{u} \tag{72}$$

with \mathbf{A}_* given by (7).

Proof Compared to the case in \mathbb{R}^m , although $\mathcal{S}_{\tilde{\mathbf{P}}} = \{\tilde{\mathbf{u}} \in \mathcal{H}^n : \text{Gr}(\tilde{\mathbf{u}}, \tilde{\mathbf{u}}) = \tilde{\mathbf{P}}\}$ is still closed and bounded ($\|\tilde{\mathbf{u}}\|_{\mathcal{H}^n}^2 \equiv \text{tr}(\tilde{\mathbf{P}})$ on $\mathcal{S}_{\tilde{\mathbf{P}}}$), $\mathcal{S}_{\tilde{\mathbf{P}}}$ is no longer compact due to the infinite dimensionality. Therefore, although the objective function is continuous and bounded from below, its infimum may not be attainable. In the following, we first assume that a global minimizer $\tilde{\mathbf{u}}_*$ exists and derive an analytic expression for it. After that, we will show that this expression, which is well-defined whether or not a global minimum exists, is indeed the unique global minimizer to (70).

Here the objective function is

$$F(\tilde{\mathbf{u}}) = \|\tilde{\mathbf{u}} - \mathbf{u}\|_{\mathcal{H}^n}^2 = \text{tr}(\text{Gr}(\tilde{\mathbf{u}} - \mathbf{u}, \tilde{\mathbf{u}} - \mathbf{u})) = \text{tr}(\tilde{\mathbf{P}}) + \text{tr}(\mathbf{P}) - 2\text{tr}(\text{Gr}(\tilde{\mathbf{u}}, \mathbf{u})),$$

so its first-order variation is

$$\delta F(\tilde{\mathbf{u}}, \delta \tilde{\mathbf{u}}) = -2\text{tr}(\text{Gr}(\delta \tilde{\mathbf{u}}, \mathbf{u})).$$

The constraint on $\tilde{\mathbf{u}}$ is $\text{Gr}(\tilde{\mathbf{u}}, \tilde{\mathbf{u}}) = \tilde{\mathbf{P}}$, so the constraint on $\delta \tilde{\mathbf{u}}$ is

$$\text{Gr}(\delta \tilde{\mathbf{u}}, \tilde{\mathbf{u}}) + \text{Gr}(\tilde{\mathbf{u}}, \delta \tilde{\mathbf{u}}) = 0, \tag{73}$$

which is simply saying that $\text{Gr}(\delta \tilde{\mathbf{u}}, \tilde{\mathbf{u}})$ is anti-symmetric. If $\tilde{\mathbf{u}}$ is a local minimizer to (72), we must have $\delta F(\tilde{\mathbf{u}}, \delta \tilde{\mathbf{u}}) \equiv 0$ for any $\delta \tilde{\mathbf{u}} \in \mathcal{H}^n$ that satisfies (73).

Now we prove by contradiction that each entry of \mathbf{u} must be a linear combination of the entries of a local minimizer $\tilde{\mathbf{u}}$, which means $\mathbf{u} = \mathbf{A}\tilde{\mathbf{u}}$ for some $\mathbf{A} \in \mathcal{M}_{n \times n}$. If this was not the case, without loss of generality, we can assume $u_1 \notin \text{span}\{\tilde{u}_1, \dots, \tilde{u}_n\} \subset \mathcal{H}$. Therefore, the projection of u_1 onto $\text{span}\{\tilde{u}_1, \dots, \tilde{u}_n\}^\perp$ must be nonzero and we denote it by v . Now we construct $\delta\tilde{\mathbf{u}} = [v, 0, \dots, 0]^T$. We can see that $\text{Gr}(\delta\tilde{\mathbf{u}}, \tilde{\mathbf{u}}) = 0$ because $\langle v, \tilde{u}_i \rangle_{\mathcal{H}} = 0, i = 1, \dots, n$. Hence, (73) is satisfied. On the other hand, since $\langle v, u_1 \rangle_{\mathcal{H}} = \|v\|_{\mathcal{H}}^2 > 0$, we have

$$\delta F(\tilde{\mathbf{u}}, \delta\tilde{\mathbf{u}}) = -2\text{tr}(\text{Gr}(\delta\tilde{\mathbf{u}}, \mathbf{u})) = -2 \sum_{i=1}^n \langle \delta\tilde{u}_i, u_i \rangle_{\mathcal{H}} = -2 \langle v, u_1 \rangle_{\mathcal{H}} < 0,$$

which contradicts the assumption that $\tilde{\mathbf{u}}$ is a local minimizer. Therefore, any local minimizer $\tilde{\mathbf{u}}$ must satisfy $\mathbf{u} = \mathbf{A}\tilde{\mathbf{u}}$ for some $\mathbf{A} \in \mathcal{M}_{n \times n}$. Since the entries of \mathbf{u} are linearly independent, \mathbf{A} must be invertible and thus $\tilde{\mathbf{u}} = \mathbf{A}^{-1}\mathbf{u}$.

Since we have already shown that with the candidate set restricted to $\{\tilde{\mathbf{u}} = \mathbf{A}\mathbf{u}\}$, Corollary 2.2 applies. This completes the proof of (72) when a global minimum exists.

Note that no matter whether a global minimum exists or not, we can always construct $\tilde{\mathbf{u}}_* = \mathbf{A}_*\mathbf{u} \in \mathcal{S}_{\tilde{\mathbf{P}}}$ with \mathbf{A}_* given by (7). Therefore, if we can show $\|\tilde{\mathbf{u}}_* - \mathbf{u}\|_{\mathcal{H}^n}$ is indeed a lower bound of $\|\tilde{\mathbf{u}} - \mathbf{u}\|_{\mathcal{H}^n}$ on $\mathcal{S}_{\tilde{\mathbf{P}}}$, we can prove the existence and uniqueness of a global minimizer.

Since $F(\tilde{\mathbf{u}}) - \text{tr}(\tilde{\mathbf{P}}) - \text{tr}(\mathbf{P}) = -2\text{tr}(\text{Gr}(\tilde{\mathbf{u}}, \mathbf{u})) = -2 \langle \tilde{\mathbf{u}}, \mathbf{u} \rangle_{\mathcal{H}^n}$, we want to show $\langle \tilde{\mathbf{u}}, \mathbf{u} \rangle_{\mathcal{H}^n} \leq \langle \tilde{\mathbf{u}}_*, \mathbf{u} \rangle_{\mathcal{H}^n}$ for any $\tilde{\mathbf{u}} \in \mathcal{S}_{\tilde{\mathbf{P}}}$ by the Cauchy-Schwarz inequality. Since $\tilde{\mathbf{u}}_*$ is not parallel to \mathbf{u} in general, we need two linear transforms $\tilde{\mathbf{u}} = \mathbf{A}_*^{1/2}\tilde{\mathbf{v}}$ and $\mathbf{u} = \mathbf{A}_*^{-1/2}\mathbf{v}$ to make sure that the Cauchy-Schwarz inequality attains its equal sign when $\tilde{\mathbf{u}} = \mathbf{A}_*\mathbf{u}$. Therefore, we have

$$\langle \tilde{\mathbf{u}}, \mathbf{u} \rangle_{\mathcal{H}^n} = \text{tr} \left(\mathbf{A}_*^{1/2} \left\langle \tilde{\mathbf{v}}, \mathbf{v}^T \right\rangle_{\mathcal{H}} \mathbf{A}_*^{-1/2} \right) = \text{tr} \left(\left\langle \tilde{\mathbf{v}}, \mathbf{v}^T \right\rangle_{\mathcal{H}} \right) = \langle \tilde{\mathbf{v}}, \mathbf{v} \rangle_{\mathcal{H}^n} \leq \|\tilde{\mathbf{v}}\|_{\mathcal{H}^n} \|\mathbf{v}\|_{\mathcal{H}^n},$$

where the last step is the Cauchy-Schwarz inequality. On the other hand, we have

$$\begin{aligned} \|\mathbf{v}\|_{\mathcal{H}^n}^2 &= \left\langle \mathbf{A}_*^{1/2}\mathbf{u}, \mathbf{A}_*^{1/2}\mathbf{u} \right\rangle_{\mathcal{H}^n} = \text{tr} \left(\mathbf{A}_*^{1/2} \left\langle \tilde{\mathbf{u}}, \mathbf{u}^T \right\rangle_{\mathcal{H}} \mathbf{A}_*^{1/2} \right) = \text{tr}(\mathbf{A}_*\mathbf{P}), \\ \|\tilde{\mathbf{v}}\|_{\mathcal{H}^n}^2 &= \left\langle \mathbf{A}_*^{-1/2}\tilde{\mathbf{u}}, \mathbf{A}_*^{-1/2}\tilde{\mathbf{u}} \right\rangle_{\mathcal{H}^n} = \text{tr} \left(\tilde{\mathbf{P}}\mathbf{A}_*^{-1} \right) = \text{tr}(\mathbf{A}_*\mathbf{P}), \\ \langle \tilde{\mathbf{u}}_*, \mathbf{u} \rangle_{\mathcal{H}^n} &= \langle \mathbf{A}_*\mathbf{u}, \mathbf{u} \rangle_{\mathcal{H}^n} = \text{tr}(\mathbf{A}_*\mathbf{P}), \end{aligned}$$

where the last equal sign in the second line is due to the second-moment matching constraint $\mathbf{A}_*\mathbf{P}\mathbf{A}_* = \tilde{\mathbf{P}}$. Combining the above results, we have shown

$$\langle \tilde{\mathbf{u}}, \mathbf{u} \rangle_{\mathcal{H}^n} \leq \text{tr}(\mathbf{A}_*\mathbf{P}) = \langle \tilde{\mathbf{u}}_*, \mathbf{u} \rangle_{\mathcal{H}^n}$$

for any $\tilde{\mathbf{u}} \in \mathcal{S}_{\tilde{\mathbf{P}}}$, which completes the proof. □

Theorem 2.2 can also be generalized to the case with \mathbb{R}^m replaced by a Hilbert space \mathcal{H} . This continuous setup can be intuitively viewed as the discrete one with

$m = \infty$, which implies that we always have $m > n \geq \tilde{r}$ so the modification in step 1 is never needed. We state this generalization in the following theorem, whose proof is omitted since it is totally in parallel with that of Theorem 2.2.

Theorem B.2 (Degenerate counterpart of Theorem B.1) *Given $\mathbf{u} \in \mathcal{H}^n$ with $\mathbf{P} = \text{Gr}(\mathbf{u}, \mathbf{u}) \in \mathcal{M}_{n \times n}$ and $\text{rank}(\mathbf{P}) = \dim(\text{span}\{u_1, \dots, u_n\}) = r$, and semi-SPD $\tilde{\mathbf{P}} \in \mathcal{M}_{n \times n}$ with $\text{rank}(\tilde{\mathbf{P}}) = \tilde{r}$, we have*

$$\arg \min_{\tilde{\mathbf{u}} \in \mathcal{H}^n: \text{Gr}(\tilde{\mathbf{u}}, \tilde{\mathbf{u}}) = \tilde{\mathbf{P}}} \|\tilde{\mathbf{u}} - \mathbf{u}\|_{\mathcal{H}^n}^2 = \tilde{\mathbf{u}}_* = \tilde{\mathbf{Z}}\tilde{\mathbf{s}}_* = \tilde{\mathbf{Z}}(\mathbf{Z}\tilde{\mathbf{w}}_* + \mathbf{Z}_\perp\tilde{\mathbf{w}}_{\perp*}) \tag{74}$$

where $\tilde{\mathbf{s}}_* = \mathbf{Z}\tilde{\mathbf{w}}_* + \mathbf{Z}_\perp\tilde{\mathbf{w}}_{\perp*} \in \mathcal{H}^{\tilde{r}}$. Moreover, the columns of $\tilde{\mathbf{Z}} \in \mathcal{M}_{n \times \tilde{r}}$ form an orthonormal basis of $\text{Row}(\tilde{\mathbf{P}})$ and the columns of $\mathbf{Z} \in \mathcal{M}_{\tilde{r} \times r'}$ form an orthonormal basis of $\text{Row}(\tilde{\mathbf{Z}}^T \mathbf{P} \tilde{\mathbf{Z}})$ with

$$r' = \tilde{r} - \dim(\text{Row}(\tilde{\mathbf{P}}) \cap \text{Row}(\mathbf{P})^\perp) = r - \dim(\text{Row}(\mathbf{P}) \cap \text{Row}(\tilde{\mathbf{P}})^\perp). \tag{75}$$

Finally, $\tilde{\mathbf{w}}_*$ and $\tilde{\mathbf{w}}_{\perp*}$ are given by

$$\tilde{\mathbf{w}}_* = \arg \min_{\tilde{\mathbf{w}} \in \mathcal{H}^{r'}: \text{Gr}(\tilde{\mathbf{w}}, \tilde{\mathbf{w}}) = \mathbf{Z}^T \tilde{\mathbf{Z}}^T \tilde{\mathbf{P}} \tilde{\mathbf{Z}} \mathbf{Z}} \|\tilde{\mathbf{w}} - \mathbf{Z}^T \tilde{\mathbf{Z}}^T \mathbf{u}\|_{\mathcal{H}^{r'}}^2 \tag{76}$$

and $\tilde{\mathbf{w}}_{\perp*} = \mathbf{B}_{12}^T \mathbf{B}_{11}^{-1} \tilde{\mathbf{w}}_* + \mathbf{v}_* \in \mathcal{H}^{\tilde{r}-r'}$ with arbitrary \mathbf{v}_* such that $\text{Gr}(\mathbf{v}_*, \tilde{\mathbf{w}}_*) = \mathbf{0}$, i.e. $v_{*,i} \in \text{span}\{\tilde{w}_{*,1}, \dots, \tilde{w}_{*,r'}\}^\perp \subset \mathcal{H}$ for $i = 1, \dots, (\tilde{r} - r')$, and $\text{Gr}(\mathbf{v}_*, \mathbf{v}_*) = \mathbf{B}_{22} - \mathbf{B}_{12}^T \mathbf{B}_{11}^{-1} \mathbf{B}_{12}$. The \mathbf{B}_{ij} blocks are defined as:

$$\begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{12}^T & \mathbf{B}_{22} \end{bmatrix} \triangleq \text{Gr} \left(\begin{bmatrix} \tilde{\mathbf{w}} \\ \tilde{\mathbf{w}}_{\perp} \end{bmatrix}, \begin{bmatrix} \tilde{\mathbf{w}} \\ \tilde{\mathbf{w}}_{\perp} \end{bmatrix} \right) = \begin{bmatrix} \mathbf{Z}^T \\ \mathbf{Z}_\perp^T \end{bmatrix} \tilde{\mathbf{Z}}^T \tilde{\mathbf{P}} \tilde{\mathbf{Z}} \begin{bmatrix} \mathbf{Z} & \mathbf{Z}_\perp \end{bmatrix}.$$

Besides, the global minimum of the objective function is

$$\|\tilde{\mathbf{u}}_* - \mathbf{u}\|_{\mathcal{H}^n}^2 = \text{tr}[\mathbf{P} + \tilde{\mathbf{P}} - 2(\mathbf{P}\tilde{\mathbf{P}})^{1/2}]. \tag{77}$$

All previous remarks regarding Theorem 2.2 apply here as well. Moreover, this theorem also extends our optimal transport interpretation in ‘‘Appendix C’’ to the most general case with arbitrary semi-SPD covariance matrix Σ_μ and Σ_ν , unlike most previous literature (e.g. [12,40]) on this topic, which essentially restrict themselves in cases with $\text{Row}(\Sigma_\nu) \subset \text{Row}(\Sigma_\mu)$. In particular, Theorem B.2 implies that the optimal coupling π_* between μ and ν (which is a joint distribution with μ and ν being its marginals) for matching only the first two moments is deterministic (i.e. corresponds to a bijection between $\mathbf{u} \sim \mu$ and $\tilde{\mathbf{u}} \sim \nu$) if and only if $\text{rank}(\Sigma_\mu) = \text{rank}(\Sigma_\nu) = \text{rank}(\Sigma_\mu \Sigma_\nu)$. The optimal coupling is deterministic in the direction from μ to ν (i.e. π_* induces a map from \mathbf{u} to $\tilde{\mathbf{u}}$ or the conditional $\pi_{\nu|\mu}$ is always a Dirac delta) if and only if $\text{Row}(\Sigma_\nu) \cap \text{Row}(\Sigma_\mu)^\perp = \mathbf{0}$, of which the assumption $\text{Row}(\Sigma_\nu) \subset \text{Row}(\Sigma_\mu)$ usually made in previous literature is a special case.

C: Connections to optimal transport

Our main results Theorems 2.1 and B.1 are also closely connected to the problem of optimal transport (a.k.a. the Monge-Kantorovich minimization problem [55, p. 10]), which concerns the coupling (in the form of a joint distribution) between two given probability distributions that minimizes an associated transport cost. More precisely, in a particular setting that is relevant to our context, given two PDFs $\mu(\cdot)$ and $\nu(\cdot)$ with finite second moments over \mathbb{R}^n , the problem concerns the joint distribution $\pi(\cdot, \cdot)$ over $\mathbb{R}^n \times \mathbb{R}^n$, which can be degenerate, such that π has μ and ν as marginals and the L_2 distance between the two marginal random variables are minimized, i.e.

$$\pi_* = \arg \min_{\pi \in \Pi_{\mu, \nu}} \int_{\mathbb{R}^n \times \mathbb{R}^n} \|x - y\|^2 \pi(x, y) dx dy, \tag{78}$$

where $\|\cdot\|$ is the Euclidean distance in \mathbb{R}^n and

$$\Pi_{\mu, \nu} = \left\{ \pi(\cdot, \cdot) \in \text{PDFs} : \int_{\mathbb{R}^n} \pi(x, y) dy = \mu(x), \int_{\mathbb{R}^n} \pi(x, y) dx = \nu(y) \right\}. \tag{79}$$

Under the above conditions, a unique optimizer π_* exists [55, p. 11] and it turns out to be deterministic, i.e. $\pi_*(x, y) = \delta_{y=T_*(x)}\mu(x)$ for some mapping $T_* : \mathbb{R}^n \rightarrow \mathbb{R}^n$ such that $\int \pi_*(x, y) dx = \nu(y)$. Moreover, the quantity

$$W_2(\mu, \nu) = \min_{\pi \in \Pi_{\mu, \nu}} \left(\int_{\mathbb{R}^n \times \mathbb{R}^n} \|x - y\|^2 \pi(x, y) dx dy \right)^{1/2} \tag{80}$$

is called the L_2 Wasserstein distance between μ and ν [55, ch. 6].

In general, the optimal transport mapping $T_*(\cdot)$ is nonlinear. However, if we relax the constraint of marginal matching to only first- and second-moment matching, Theorem B.1 implies that $T_*(\cdot)$ is indeed affine with the linear part being SPD, even when μ and ν do not have a well-defined PDF (i.e. not absolutely continuous with respect to the Lebesgue measure). To show this, first notice that the optimization (78) over joint distributions on $\mathbb{R}^n \times \mathbb{R}^n$ can be translated into an equivalent problem over random vectors in \mathbb{R}^n :

$$\tilde{u}_* = \arg \min_{\tilde{u} \sim \nu} \mathbb{E} \left[\|\tilde{u} - u\|^2 \right] \tag{81}$$

with some $u \sim \mu$. Next, denote the mean of μ and ν by m_μ and m_ν , respectively, and also the covariance matrices by Σ_μ and Σ_ν . Then the objective function in (81) becomes

$$\mathbb{E} \left[\|\tilde{u} - u\|^2 \right] = \|m_\mu - m_\nu\|^2 + \mathbb{E} \left[\|\tilde{u}' - u'\|^2 \right],$$

where the primes indicate the centered random vectors. Now if we relax the distribution matching constraint $\tilde{u} \sim \nu$ to only mean and covariance matching, i.e. $\mathbb{E}[\tilde{u}] = m_\nu$

and $\mathbb{E}[\tilde{\mathbf{u}}' \tilde{\mathbf{u}}'^T] = \Sigma_v$, the optimal transport problem (81) is simplified to the task of minimum-correction second-moment matching:

$$\tilde{\mathbf{u}}'_* = \arg \min_{\mathbb{E}[\tilde{\mathbf{u}}' \tilde{\mathbf{u}}'^T] = \Sigma_v} \mathbb{E} \left[\|\tilde{\mathbf{u}} - \mathbf{u}\|^2 \right]. \tag{82}$$

If we set $\mathcal{H} = L_2(\Omega)$ as the Hilbert space of all square-integrable random variables defined over the probability space $(\Omega, \mathcal{A}, \mathcal{P})$ on which \mathbf{u} is defined, then Theorem B.1 applies to (82) and yields $\tilde{\mathbf{u}}_* = \mathbf{T}_*(\mathbf{u}) = \mathbf{A}_*(\mathbf{u} - \mathbf{m}_\mu) + \mathbf{m}_v$ with

$$\mathbf{A}_* = \sqrt{\Sigma_\mu}^{-1} (\sqrt{\Sigma_\mu} \Sigma_v \sqrt{\Sigma_\mu}^T)^{1/2} \sqrt{\Sigma_\mu}^{-T}. \tag{83}$$

Moreover, since

$$\mathbb{E} \left[\|\tilde{\mathbf{u}}' - \mathbf{u}'\|^2 \right] = \text{tr}(\Sigma_\mu) + \text{tr}(\Sigma_v) - 2\text{tr}[\text{Cov}(\tilde{\mathbf{u}}, \mathbf{u})] \tag{84}$$

and

$$\text{tr}[\text{Cov}(\tilde{\mathbf{u}}_*, \mathbf{u})] = \text{tr}[\text{Cov}(\mathbf{A}_* \mathbf{u}', \mathbf{u}')] = \text{tr}[(\sqrt{\Sigma_\mu} \Sigma_v \sqrt{\Sigma_\mu}^T)^{1/2}] = \text{tr}[(\Sigma_\mu \Sigma_v)^{1/2}],$$

we have

$$W_2^2(\mu, \mathbf{T}_* \mu) = \|\mathbf{m}_\mu - \mathbf{m}_v\|^2 + \text{tr}[\Sigma_\mu + \Sigma_v - 2(\Sigma_\mu \Sigma_v)^{1/2}], \tag{85}$$

where the distribution $\mathbf{T}_* \mu$ is the push-forward of μ by \mathbf{T}_* . Note that although $\mathbf{T}_* \mu$ share the same mean and covariance with v , in general $\mathbf{T}_* \mu \neq v$ and $W_2(\mu, \mathbf{T}_* \mu) < W_2(\mu, v)$, so (85) indeed provides a lower bound for $W_2(\mu, v)$ based on the first two moments of μ and v .

If we partition the set of all distributions (possibly degenerate) over \mathbb{R}^n that have finite second moments into equivalent classes by their means and covariances then (85) induces a metric on this quotient space. Denote by $\mathcal{S}_{\mathbf{m}, \Sigma}$ the equivalent class of all distributions with mean \mathbf{m} and covariance Σ . Given $\mathcal{S}_{\mathbf{m}_2, \Sigma_2}$ and $\mu \in \mathcal{S}_{\mathbf{m}_1, \Sigma_1}$, there exists a unique $\tilde{\mu} = \mathbf{T}_* \mu \in \mathcal{S}_{\mathbf{m}_2, \Sigma_2}$ that is the closest to μ under the distance $W_2(\cdot, \cdot)$. As a generalization to the one-to-one correspondence of (9) and the intuition illustrated in Fig. 1, the members in any two equivalent classes $\mathcal{S}_{\mathbf{m}_1, \Sigma_1}$ and $\mathcal{S}_{\mathbf{m}_2, \Sigma_2}$ can be paired up by such minimum- W_2 -distance matching. Moreover, since this distance between μ and $\mathbf{T}_* \mu$ depends only on $\mathcal{S}_{\mathbf{m}_1, \Sigma_1}$ and $\mathcal{S}_{\mathbf{m}_2, \Sigma_2}$ but not on the particular member μ , $W_2(\mu, \mathbf{T}_* \mu)$ can be used to define the distance between $\mathcal{S}_{\mathbf{m}_1, \Sigma_1}$ and $\mathcal{S}_{\mathbf{m}_2, \Sigma_2}$. This is sometimes called the Fréchet distance [12]. Note that the covariance part of (85) can also be isolated and used as a metric among SPD matrices.

Although in general (85) is only a lower bound for $W_2(\mu, v)$, if both μ and v belong to a family of distributions closed under affine transforms and each member of this family is uniquely determined by its mean and covariance, then we have $\mathbf{T}_* \mu = v$ and thus $W_2(\mu, v)$ coincides with (85). [12] mentions one typical example of such where the distribution family is characterized by a radial kernel $\phi : [0, \infty) \rightarrow [0, \infty)$

that satisfies $0 < \int_0^\infty r^{2n+1} \phi(r^2) dr < \infty$, which guarantees the finiteness of second moments. In this family, a member μ 's PDF takes the form $\mu(\mathbf{x}) \propto \phi((\mathbf{x} - \mathbf{m})^T \mathbf{A}(\mathbf{x} - \mathbf{m}))$ with some SPD \mathbf{A} , where \mathbf{m} is the mean and \mathbf{A}^{-1} is a multiple of the covariance. In other words, this family is generated by pushing forward a non-degenerate elliptically contoured distribution that has finite second moments by invertible affine transforms. It can be checked that within such a family, a distribution can be uniquely identified by its mean and covariance. Therefore, given μ and ν , we always have $\mathbf{T}_* \mu = \nu$ and thus $W_2(\mu, \nu)$ can be computed by (85). The most familiar example of such is probably the family of Gaussian distributions corresponding to the exponential kernel $\phi(r) = e^{-r}$, for which the above results are well known.

To the best of our knowledge, the earliest work that derives the right hand side expression in (85) is in [12,40]. However, they did not connect this result to the problem of optimal transport. Moreover, they translate the optimization into one over the cross-covariance matrix $\mathbf{B} = \text{Cov}(\mathbf{u}, \tilde{\mathbf{u}})$, which does not necessarily uniquely determine $\tilde{\mathbf{u}}$ given \mathbf{u} if the optimizer \mathbf{B}_* does not correspond to a deterministic affine coupling between the two. In Theorem B.1, although the objective function is also simplified to $\text{tr}(\text{Cov}(\mathbf{u}, \tilde{\mathbf{u}}))$ (or its empirical counterpart (59) in the discrete case of Theorem 2.1), we have provided an alternative proof based on optimizing on $\tilde{\mathbf{u}}$ directly rather than on $\mathbf{B} = \text{Cov}(\mathbf{u}, \tilde{\mathbf{u}})$.

Finally, note also that the more general result Theorem B.1 on a Hilbert space indeed contains the special case on \mathbb{R}^m described by Theorem 2.1. Moreover, if we view Theorem 2.1 through the lens of Theorem B.1, we can identify a discrete counterpart of the optimal transport interpretation for Theorem 2.1. The key observation is that Theorem 2.1 is equivalent to Theorem B.1 with $\mathcal{H} = \mathbb{R}^m$, which on the other hand, corresponds to (82) with μ being a discrete distribution over \mathbb{R}^n with m particles, i.e. with its (generalized) PDF being the sum of m Dirac deltas. Therefore, the optimal transport interpretation of Theorem B.1 implies that $\tilde{\mathbf{u}}_*$ also follows a discrete distribution with its m particles given by the rows of $\tilde{\mathbf{U}}_*$ in Theorem 2.1. This is actually a stronger result than Theorem 2.1 because in the latter, we have restricted the feasible $\tilde{\mu}$'s within the m -particle discrete distributions a priori.

For an exposition of other connections between the optimal transport and the polar decomposition, see [4,7].

D: Gradient descent for matrix square root

[18] proposes to solve the equation $\mathbf{X}^T \mathbf{P} \mathbf{X} = \mathbf{I}$ for an \mathbf{X} close to \mathbf{I} by solving the unconstrained optimization $\arg \min_{\mathbf{X}} \|\mathbf{X}^T \mathbf{P} \mathbf{X} - \mathbf{I}\|_F^2$ iteratively with the initial guess $\mathbf{X}_0 = \mathbf{I}$.

The objective function attains its global minimum 0 at every \mathbf{X} that solves $\mathbf{X}^T \mathbf{P} \mathbf{X} = \mathbf{I}$. The hope is that if we use an iterative scheme and start from \mathbf{I} , the iterates will converge to a solution \mathbf{X} close to \mathbf{I} . [18] proposes to use a simple gradient descent iteration. Since $f(\mathbf{X}) = \|\mathbf{X}^T \mathbf{P} \mathbf{X} - \mathbf{I}\|^2 = \text{tr}(\mathbf{X}^T \mathbf{P} \mathbf{X} \mathbf{X}^T \mathbf{P} \mathbf{X}) - 2\text{tr}(\mathbf{X}^T \mathbf{P} \mathbf{X}) + \text{tr}(\mathbf{I})$,

we can compute the gradient analytically as

$$\nabla_X f = 4\mathbf{P}\mathbf{X}(\mathbf{X}^T\mathbf{P}\mathbf{X} - \mathbf{I}). \tag{86}$$

Therefore, the gradient descent iteration takes the form

$$\mathbf{X}_{n+1} = \mathbf{X}_n - \gamma\mathbf{P}\mathbf{X}_n(\mathbf{X}_n^T\mathbf{P}\mathbf{X}_n - \mathbf{I}) \tag{87}$$

with γ being a tunable step size. This is Algorithm 3 in [18, sec. 3.5]. However, [18] does not mention how to pick γ , when the iteration converges, and how fast.

We can analyze this iteration by exactly the same techniques used in Sect. 3.2. First, we have $\mathbf{X}_k \rightarrow \mathbf{P}^{-1/2}\mathbf{Q}_0$ for

$$\mathbf{X}_0 = \mathbf{S}_0\mathbf{Q}_0, \quad \mathbf{S}_0 \in \mathcal{C}_\mathbf{P}, \quad \|\mathbf{X}_0\|_2 \leq 1/\|\mathbf{P}\|_2^{1/2}, \quad \gamma \leq 1/(2\|\mathbf{P}\|_2), \tag{88}$$

since the matrix iteration can be reduced to the scalar one $\sigma_{k+1} = \sigma_k - \gamma\lambda\sigma_k(\sigma_k^2 - 1)^8$, where λ is an eigenvalue of \mathbf{P} and $\sigma_k/\sqrt{\lambda}$ is the corresponding singular value of \mathbf{X}_k . Next, the stability analysis shows that the Fréchet derivative at the limit \mathbf{X}_* is

$$\mathbf{D}\mathbf{F}(\mathbf{X}) = \mathbf{X} - \gamma\mathbf{A}\mathbf{X} - \gamma\mathbf{A}^{1/2}\mathbf{X}^T\mathbf{A}^{1/2}, \tag{89}$$

which has eigenvectors and eigenvalues as

$$\begin{aligned} \mathbf{Z}_{ij} &= \mathbf{E}_{ij} + \sqrt{\lambda_j/\lambda_i}\mathbf{E}_{ji}, & \mu_{ij} &= 1 - \gamma(\lambda_i + \lambda_j), & i \leq j, \\ \mathbf{Z}_{ij} &= \mathbf{E}_{ij} - \sqrt{\lambda_i/\lambda_j}\mathbf{E}_{ji}, & \mu_{ij} &= 1, & i > j. \end{aligned} \tag{90}$$

Here the λ_i 's are the diagonal entries of \mathbf{A} and are also the eigenvalues of \mathbf{P} . Therefore, for $\rho(\mathbf{D}\mathbf{F}) < 1$, we need $\gamma < 1/\|\mathbf{P}\|_2$, which is less stringent than (88). Moreover, the iteration converges linearly with the asymptotic error diminishing coefficient not smaller than $(\lambda_{\max} - \lambda_{\min})/(\lambda_{\max} + \lambda_{\min})$. Therefore, the gradient descent is considered unattractive here compared to the quadratically converging iterations (45) and (46).

References

1. Absil, P.A., Mahony, R., Sepulchre, R.: Optimization Algorithms on Matrix Manifolds. Princeton University Press, Princeton (2009)
2. Absil, P.A., Malick, J.: Projection-like retractions on matrix manifolds. SIAM J. Optim. **22**(1), 135–158 (2012)
3. Anderson, J.L.: An ensemble adjustment Kalman filter for data assimilation. Mon. Weather Rev. **129**(12), 2884–2903 (2001)
4. Benamou, J.D., Brenier, Y.: A computational fluid mechanics solution to the Monge-Kantorovich mass transfer problem. Numerische Mathematik **84**(3), 375–393 (2000)
5. Bergemann, K., Gottwald, G., Reich, S.: Ensemble propagation and continuous matrix factorization algorithms. Q. J. R. Meteorol. Soc. **135**(643), 1560–1572 (2009)

⁸ Since $(1 - \sigma_{k+1}) = (1 - \sigma_k)(1 - \gamma\lambda\sigma_k(\sigma_k + 1))$, if $0 < \sigma_0 \leq 1$ and $0 < \gamma \leq 1/(2\lambda)$, σ_k monotonically increases and converges to 1.

6. Bishop, C.H., Etherton, B.J., Majumdar, S.J.: Adaptive sampling with the ensemble transform Kalman filter. Part I: theoretical aspects. *Mon. Weather Rev.* **129**(3), 420–436 (2001)
7. Brenier, Y.: Polar factorization and monotone rearrangement of vector-valued functions. *Commun. Pure Appl. Math.* **44**(4), 375–417 (1991)
8. Brockett, R.W.: Dynamical systems that learn subspaces. In: *Mathematical System Theory*, pp. 579–592. Springer (1991)
9. Chandrasekaran, S., Ipsen, I.C.: Backward errors for eigenvalue and singular value decompositions. *Numerische Mathematik* **68**(2), 215–223 (1994)
10. Cheng, M., Hou, T.Y., Zhang, Z.: A dynamically bi-orthogonal method for time-dependent stochastic partial differential equations I: derivation and algorithms. *J. Comput. Phys.* **242**, 843–868 (2013)
11. Dieci, L., Russell, R.D., Van Vleck, E.S.: Unitary integrators and applications to continuous orthonormalization techniques. *SIAM J. Numer. Anal.* **31**(1), 261–281 (1994)
12. Dowson, D., Landau, B.: The Fréchet distance between multivariate normal distributions. *J. Multivar. Anal.* **12**(3), 450–455 (1982)
13. Dutt, A.: High order stochastic transport and Lagrangian data assimilation. Master's thesis, Massachusetts Institute of Technology, Department of Mechanical Engineering, Cambridge, Massachusetts (2018)
14. Edelman, A., Arias, T.A., Smith, S.T.: The geometry of algorithms with orthogonality constraints. *SIAM J. Matrix Anal. Appl.* **20**(2), 303–353 (1998)
15. Evensen, G.: Sequential data assimilation with a nonlinear quasi-geostrophic model using Monte Carlo methods to forecast error statistics. *J. Geophys. Res. Oceans* **99**(C5), 10143–10162 (1994)
16. Evensen, G.: The ensemble Kalman filter: theoretical formulation and practical implementation. *Ocean Dyn.* **53**(4), 343–367 (2003)
17. Feppon, F., Lermusiaux, P.F.J.: A geometric approach to dynamical model-order reduction. *SIAM J. Matrix Anal. Appl.* **39**(1), 510–538 (2018). <https://doi.org/10.1137/16M1095202>
18. Feppon, F., Lermusiaux, P.F.J.: Dynamically orthogonal numerical schemes for efficient stochastic advection and Lagrangian transport. *SIAM Rev.* **60**(3), 595–625 (2018). <https://doi.org/10.1137/16M1109394>
19. Feppon, F., Lermusiaux, P.F.J.: The extrinsic geometry of dynamical systems tracking nonlinear matrix projections. *SIAM J. Matrix Anal. Appl.* **40**(2), 814–844 (2019). <https://doi.org/10.1137/18M1192780>
20. Golub, G.H., Van Loan, C.F.: *Matrix Comput*, 4th edn. JHU Press, Baltimore (2013)
21. Govaerts, W., Werner, B.: Continuous bordering of matrices and continuous matrix decompositions. *Numerische Mathematik* **70**(3), 303–310 (1995)
22. Hairer, E., Lubich, C., Wanner, G.: *Geometric Numerical Integration: Structure-Preserving Algorithms for Ordinary Differential Equations*, vol. 31. Springer, Berlin (2006)
23. Hale, N., Higham, N.J., Trefethen, L.N.: Computing A^α , $\log(A)$, and related matrix functions by contour integrals. *SIAM J. Numer. Anal.* **46**(5), 2505–2523 (2008)
24. Higham, D.J.: Runge-Kutta type methods for orthogonal integration. *Appl. Numer. Math.* **22**(1–3), 217–223 (1996)
25. Higham, D.J.: Time-stepping and preserving orthonormality. *BIT Numer. Math.* **37**(1), 24–36 (1997)
26. Higham, N.J.: Computing the polar decomposition-with applications. *SIAM J. Sci. Stat. Comput.* **7**(4), 1160–1174 (1986)
27. Higham, N.J.: Newton's method for the matrix square root. *Math. Comput.* **46**(174), 537–549 (1986)
28. Higham, N.J.: *Matrix Nearness Problems and Applications*. Department of Mathematics, University of Manchester, Tech. rep (1988)
29. Higham, N.J., Schreiber, R.S.: Fast polar decomposition of an arbitrary matrix. *SIAM J. Sci. Stat. Comput.* **11**(4), 648–655 (1990)
30. Koch, O., Lubich, C.: Dynamical low-rank approximation. *SIAM J. Matrix Anal. Appl.* **29**(2), 434–454 (2007)
31. Laub, A.J.: *Matrix Analysis for Scientists and Engineers*. SIAM, New Delhi (2005)
32. Lermusiaux, P.F.J.: Error subspace data assimilation methods for ocean field estimation: theory, validation and applications. Ph.D. thesis, Harvard University, Cambridge, Massachusetts (1997)
33. Lermusiaux, P.F.J.: Data assimilation via error subspace statistical estimation, part II: Mid-Atlantic Bight shelfbreak front simulations, and ESSE validation. *Mon. Weather Rev.* **127**(7), 1408–1432 (1999). <https://doi.org/10.1175/1520-0493>
34. Lermusiaux, P.F.J.: Evolving the subspace of the three-dimensional multiscale ocean variability: Massachusetts Bay. *J. Mar. Syst.* **29**(1), 385–422 (2001)

35. Lermusiaux, P.F.J., Robinson, A.R.: Data assimilation via error subspace statistical estimation, part I: theory and schemes. *Mon. Weather Rev.* **127**(7), 1385–1407 (1999). <https://doi.org/10.1175/1520-0493>
36. Lin, J.: Bayesian learning for high-dimensional nonlinear dynamical systems: methodologies, numerics and applications to fluid flows. Ph.D. thesis, Massachusetts Institute of Technology, Department of Mechanical Engineering, Cambridge, Massachusetts (2020)
37. Lin, J.: Minimum-correction second-moment matching: theory, algorithms and applications. Master's thesis, Massachusetts Institute of Technology, Department of Mechanical Engineering, Cambridge, Massachusetts (2020)
38. Lorenz, E.N.: Predictability: a problem partly solved. In: *Proc. Seminar on Predictability*, vol. 1 (1996)
39. Musharbash, E., Nobile, F.: Dual dynamically orthogonal approximation of incompressible Navier Stokes equations with random boundary conditions. *J. Comput. Phys.* **354**, 135–162 (2018)
40. Olkin, I., Pukelsheim, F.: The distance between two random vectors with given dispersion matrices. *Linear Algebra Appl.* **48**, 257–263 (1982)
41. Ott, E., Hunt, B.R., Szunyogh, I., Zimin, A.V., Kostelich, E.J., Corazza, M., Kalnay, E., Patil, D., Yorke, J.A.: A local ensemble Kalman filter for atmospheric data assimilation. *Tellus A* **56**(5), 415–428 (2004)
42. Petersen, K.B., Pedersen, M.S.: *The matrix cookbook*. Technical University of Denmark (2012)
43. Philippe, B.: An algorithm to improve nearly orthonormal sets of vectors on a vector processor. *SIAM J. Algebr. Discrete Methods* **8**(3), 396–403 (1987)
44. Reich, S., Cotter, C.: *Probabilistic Forecasting and Bayesian Data Assimilation*. Cambridge University Press, Cambridge (2015)
45. Rheinboldt, W.C.: On the computation of multi-dimensional solution manifolds of parametrized equations. *Numerische Mathematik* **53**(1–2), 165–181 (1988)
46. Sapsis, T.P., Lermusiaux, P.F.J.: Dynamically orthogonal field equations for continuous stochastic dynamical systems. *Physica D: Nonlinear Phenomena* **238**(23–24), 2347–2360 (2009). <https://doi.org/10.1016/j.physd.2009.09.017>
47. Schatzman, M.: *Numerical Analysis: A Mathematical Introduction*. Oxford University Press, Oxford (2002)
48. Strogatz, S.H.: *Nonlinear Dynamics and Chaos: With Applications to Physics, Biology, Chemistry, and Engineering*. Westview Press, Boulder (1994)
49. Subramani, D.N.: Probabilistic regional ocean predictions: stochastic fields and optimal planning. Ph.D. thesis, Massachusetts Institute of Technology, Department of Mechanical Engineering, Cambridge, Massachusetts (2018)
50. Tagare, H.D.: *Notes on Optimization on Stiefel Manifolds*. Yale University, New Haven (2011)
51. Trefethen, L.N., Bau, D.I.I.I.: *Numerical Linear Algebra*, vol. 50. SIAM, New Delhi (1997)
52. Uecker mann, M., Lermusiaux, P., Sapsis, T.: Numerical schemes and computational studies for dynamically orthogonal equations. MSEAS Report 11, Department of Mechanical Engineering, Massachusetts Institute of Technology (2011). <http://mseas.mit.edu/?p=1928>
53. Uecker mann, M.P., Lermusiaux, P.F.J., Sapsis, T.P.: Numerical schemes for dynamically orthogonal equations of stochastic fluid and ocean flows. *J. Comput. Phys.* **233**, 272–294 (2013). <https://doi.org/10.1016/j.jcp.2012.08.041>
54. Vetra-Carvalho, S., van Leeuwen, P.J., Nerger, L., Barth, A., Altaf, M.U., Brasseur, P., Kirchgessner, P., Beckers, J.M.: State-of-the-art stochastic data assimilation methods for high-dimensional non-Gaussian problems. *Tellus A: Dyn. Meteorol. Oceanogr.* **70**(1), 1–43 (2018). <https://doi.org/10.1080/16000870.2018.1445364>
55. Villani, C.: *Optimal Transport: Old and New*, vol. 338. Springer, Berlin (2008)
56. Wen, Z., Yin, W.: A feasible method for optimization with orthogonality constraints. *Math. Program.* **142**(1–2), 397–434 (2013)
57. Yang, H., Li, H.: Weighted polar decomposition and WGL partial ordering of rectangular complex matrices. *SIAM J. Matrix Anal. Appl.* **30**(2), 898–924 (2008)