# Learning with Symmetric Positive Definite Matrices via Generalized Bures-Wasserstein Geometry

Andi Han[1], Bamdev Mishra[2], Pratik Jawanpuria[2], and Junbin Gao[1]

[1] University of Sydney
{andi.han,junbin.gao}@sydney.edu.au
[2] Microsoft India
{bamdevm, pratik.jawanpuria}@microsoft.com

**Abstract.** Learning with symmetric positive definite (SPD) matrices has many applications in machine learning. Consequently, understanding the Riemannian geometry of SPD matrices has attracted much attention lately. A particular Riemannian geometry of interest is the recently proposed Bures-Wasserstein (BW) geometry which builds on the Wasserstein distance between the Gaussian densities. In this paper, we propose a novel generalization of the BW geometry, which we call the GBW geometry. The proposed generalization is parameterized by a symmetric positive definite matrix $\mathbf{M}$ such that when $\mathbf{M} = \mathbf{I}$, we recover the BW geometry. We provide a rigorous treatment to study various differential geometric notions on the proposed novel generalized geometry which makes it amenable to various machine learning applications. We also present experiments that illustrate the efficacy of the proposed GBW geometry over the BW geometry.

**Keywords:** Riemannian geometry · SPD matrices · Bures-Wasserstein.

## 1 Introduction

Symmetric positive definite (SPD) matrices play a fundamental role in various fields of machine learning, such as metric learning [31], signal processing [12], sparse coding [13,23], computer vision [20,40], and medical imaging [46,39], etc. The set of SPD matrices, denoted as $\mathbb{S}_{++}^n$, is a subset of the Euclidean space $\mathbb{R}^{n(n+1)/2}$. To measure the (dis)similarity between SPD matrices, one needs to assign a metric (an inner product structure on the tangent space) on $\mathbb{S}_{++}^n$, which yields a Riemannian manifold. Consequently, various Riemannian metrics have been studied such as the affine-invariant [7,46], Log-Euclidean [4], and Log-Cholesky [38] metrics, and those induced from symmetric divergences [50,51]. Different metrics lead to different differential structures on the SPD matrices, and therefore, picking the "right" one depends on the application at hand. Indeed, the choice of metric has profound effect on the performance of learning algorithms [43,49,21].

The Bures-Wasserstein (BW) metric and its geometry for SPD matrices have lately gained popularity, especially in machine learning applications [8,41,54] such as statistical optimal transport [8], computer graphics [48], neural sciences [19], and evolutionary biology [14], among others. It also connects to the theory of optimal transport and the $L_2$-Wasserstein distance between zero-centered Gaussian densities [8]. More recently, [21] analyzes the BW and the affine-invariant (AI) geometries in SPD learning problems and compare their advantages/disadvantages in various machine learning applications.

In this paper, we propose a natural generalization of the BW metric by scaling SPD matrices with a given parameter SPD matrix $\mathbf{M}$. The introduction of $\mathbf{M}$ gives flexibility to the BW metric. Choosing $\mathbf{M}$ is equivalent to choosing a suitable metric for learning tasks on SPD matrices. For example, a proper choice of $\mathbf{M}$ can lead to faster convergence of algorithms for certain class of optimization problems (see more discussions in Section 4). Indeed, when $\mathbf{M} = \mathbf{I}$, the generalized metric reduces to the BW metric for SPD matrices. When $\mathbf{M} = \mathbf{X}$, the proposed metric coincides locally with the AI metric, i.e., around the neighbourhood of a SPD matrix $\mathbf{X}$. The proposed generalized metric allows to connect the BW and AI metrics (locally) with different choices of $\mathbf{M}$. The following are our contributions.

- We propose a novel generalized BW (GBW) metric by generalizing the Lyapunov operation in the BW metric (Section 2). In addition, it can also be viewed as a generalized Procrustes distance and also as the Wasserstein distance with Mahalanobis cost metric for Gaussians.
- The GBW metric leads to a Riemannian geometry for SPD matrices. In Section 3.1, we derive various Riemannian operations like geodesics, exponential and logarithm maps, Levi-Civita connection. We show that they are also natural generalizations of operations with the BW geometry. Section 3.2 derives Riemannian optimization ingredients under the proposed geometry.
- In Section 4, we show the usefulness of the GBW geometry in the applications of covariance estimation and Gaussian mixture models.

## 2   Generalized Bures-Wasserstein metric

The Bures-Wasserstein (BW) distance is defined as

$$d_{\mathrm{bw}}(\mathbf{X}, \mathbf{Y}) = \sqrt{\mathrm{tr}(\mathbf{X}) + \mathrm{tr}(\mathbf{Y}) - 2\mathrm{tr}(\mathbf{X}\mathbf{Y})^{1/2}}, \qquad (1)$$

where $\mathbf{X}$ and $\mathbf{Y}$ are SPD matrices and $\mathrm{tr}(\mathbf{X})^{1/2}$ denotes the trace of the matrix square root. It has been shown in [8,41] that the BW distance (1) induces a Riemannian metric and geometry on the manifold of SPD matrices. The BW metric that leads to the distance (1) is defined as

$$g_{\mathrm{bw}}(\mathbf{U}, \mathbf{V}) = \frac{1}{2}\mathrm{tr}(\mathcal{L}_{\mathbf{X}}[\mathbf{U}]\mathbf{V}) = \frac{1}{2}\mathrm{vec}(\mathbf{U})^{\top}(\mathbf{X} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{X})^{-1}\mathrm{vec}(\mathbf{V}), \quad (2)$$

where $\mathbf{U}, \mathbf{V}$ on $T_{\mathbf{X}}\mathbb{S}_{++}^n$ are the symmetric matrices and the Lyapunov operator $\mathcal{L}_{\mathbf{X}}[\mathbf{U}]$ is defined as the solution of the matrix equation $\mathbf{X}\mathcal{L}_{\mathbf{X}}[\mathbf{U}] + \mathcal{L}_{\mathbf{X}}[\mathbf{U}]\mathbf{X} = \mathbf{U}$

for $\mathbf{U} \in \mathbb{S}^n$ (which is the set of symmetric matrices of size $n \times n$). Here, $\text{vec}(\mathbf{U})$ and $\text{vec}(\mathbf{V})$ are the vectorization of matrices $\mathbf{U}$ and $\mathbf{V}$, respectively, and $\otimes$ denotes the Kronecker product. Our proposed GBW metric generalizes (2) and is parameterized by a given $\mathbf{M} \in \mathbb{S}^n_{++}$ as

$$g_{\text{gbw}}(\mathbf{U}, \mathbf{V}) = \tfrac{1}{2}\text{tr}(\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{V}) = \tfrac{1}{2}\text{vec}(\mathbf{U})^\top (\mathbf{X} \otimes \mathbf{M} + \mathbf{M} \otimes \mathbf{X})^{-1}\text{vec}(\mathbf{V}), \quad (3)$$

where $\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]$ is the generalized Lyapunov operator, defined as the solution to the linear matrix equation $\mathbf{X}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{M} + \mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{X} = \mathbf{U}$. Similar to the special Lyapunov operator, the solution is symmetric given that $\mathbf{X}, \mathbf{M} \in \mathbb{S}^n_{++}$ and $\mathbf{U} \in \mathbb{S}^n$. As we show later that the Riemannian distance associated with the GBW metric is derived as

$$d_{\text{gbw}}(\mathbf{X}, \mathbf{Y}) = \sqrt{\text{tr}(\mathbf{M}^{-1}\mathbf{X}) + \text{tr}(\mathbf{M}^{-1}\mathbf{Y}) - 2\text{tr}(\mathbf{X}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1})^{1/2}}, \quad (4)$$

which can be seen as the BW distance (1) between $\mathbf{M}^{-1/2}\mathbf{X}\mathbf{M}^{-1/2}$ and $\mathbf{M}^{-1/2}\mathbf{Y}\mathbf{M}^{-1/2}$. Note that the affine-invariant metric [7] is given by $g_{\text{ai}}(\mathbf{U}, \mathbf{V}) = \text{vec}(\mathbf{U})^\top (\mathbf{X} \otimes \mathbf{X})^{-1}\text{vec}(\mathbf{V})$. Clearly, the proposed metric (3) coincides locally with the affine-invariant (AI) metric when $\mathbf{M} = \mathbf{X}$, i.e., around the neighbourhood of $\mathbf{X}$. (Implications of this observation are discussed later in experiments.)

Below, we show that the same GBW distance (4) is realized under various contexts naturally. In those cases, the Euclidean norm, denoted by $\|\cdot\|_2$ is replaced with the more general Mahalanobis norm defined as $\|\mathbf{X}\|_{\mathbf{M}^{-1}} := \sqrt{\text{tr}(\mathbf{X}^\top \mathbf{M}^{-1}\mathbf{X})}$.

***Orthogonal Procrustes problem:*** Any SPD matrix $\mathbf{X} \in \mathbb{S}^n_{++}$ can be factorized as $\mathbf{X} = \mathbf{P}\mathbf{P}^\top$ for $\mathbf{P} \in \text{M}(n)$, the set of invertible matrices. Such a factorization is invariant under the action of the orthogonal group $O(n)$, the set of orthogonal matrices. That is, for any $\mathbf{O} \in O(n)$, $\mathbf{P}\mathbf{O}$ is also a valid parameterization. In [8], the BW distance is verified as the extreme solution of the orthogonal Procrustes problem where $\mathbf{P}$ is set to be $\mathbf{X}^{1/2}$, i.e., $d_{\text{bw}}(\mathbf{X}, \mathbf{Y}) = \min_{\mathbf{O} \in O(n)} \|\mathbf{X}^{1/2} - \mathbf{Y}^{1/2}\mathbf{O}\|_2$. We can show that the GBW distance is obtained as the solution to the same orthogonal Procrustes problem in the Mahalanobis norm parameterized by $\mathbf{M}^{-1}$.

**Proposition 1.** $d_{\text{gbw}}(\mathbf{X}, \mathbf{Y}) = \min_{\mathbf{O} \in O(n)} \|\mathbf{X}^{1/2} - \mathbf{Y}^{1/2}\mathbf{O}\|_{\mathbf{M}^{-1}}$.

***Wasserstein distance and optimal transport:*** To demonstrate the connection of the GBW distance to the Wasserstein distance, recall that the $L_2$-Wasserstein distance between two probability measures $\mu, \nu$ with finite second moments is $W^2(\mu, \nu) = \inf_{\mathbf{x} \sim \mu, \mathbf{y} \sim \nu} \mathbb{E}\|\mathbf{x} - \mathbf{y}\|_2^2 = \inf_{\gamma \sim \Gamma(\mu,\nu)} \int_{\mathbb{R}^n \times \mathbb{R}^n} \|\mathbf{x} - \mathbf{y}\|_2^2 d\gamma(\mathbf{x}, \mathbf{y})$, where $\Gamma(\mu, \nu)$ is the set of all probability measures with marginals $\mu, \nu$. It is well known that the $L_2$-Wasserstein distance between two zero-centered Gaussian distributions is equal to the BW distance between their covariance matrices [8,47,54]. The following proposition shows that the $L_2$-Wasserstein distance between such measures with respect to a Mahalanobis cost metric (which we term as generalized Wasserstein distance) coincides with the GBW distance in (4).

**Proposition 2.** *Define the generalized Wasserstein distance as $\tilde{W}^2(\mu, \nu) :=$ $\inf_{\mathbf{x} \sim \mu, \mathbf{y} \sim \nu} \mathbb{E} \|\mathbf{x} - \mathbf{y}\|^2_{\mathbf{M}^{-1}}$, for any $\mathbf{M}^{-1} \in \mathbb{S}^n_{++}$. Suppose $\mu, \nu$ are two Gaussian measures with zero mean and covariances as $\mathbf{X}, \mathbf{Y} \in \mathbb{S}^n_{++}$ respectively. Then, we have $\tilde{W}^2(\mu, \nu) = d^2_{\mathrm{gbw}}(\mathbf{X}, \mathbf{Y})$.*

Alternatively, the same distance is recovered by considering two scaled random Gaussian vector $\mathbf{M}^{-1/2}\mathbf{x}, \mathbf{M}^{-1/2}\mathbf{y}$ under the Euclidean distance, i.e., $d^2(\mathbf{X}, \mathbf{Y}) = \inf_{\mathbf{x} \sim \mu, \mathbf{y} \sim \nu} \mathbb{E} \|\mathbf{M}^{-1/2}\mathbf{x} - \mathbf{M}^{-1/2}\mathbf{y}\|^2_2$. For completeness, we also derive the optimal transport plan corresponding to the GBW distance in Appendix.

## 3    Generalized Bures-Wasserstein Riemannian geometry

In this section, the geometry arising from the GBW metric (3) is shown to have a Riemannian structure for a given $\mathbf{M} \in \mathbb{S}^n_{++}$, which we denote as $\mathcal{M}_{\mathrm{gbw}}$. We show the expressions of the Riemannian distance, geodesic, exponential/logarithm maps, Levi-Civita connection, sectional curvature as well as the geometric mean and barycenter. A summary of the results is presented in Table 1. Additionally, we discuss optimization on the SPD manifold with the proposed GBW geometry. We defer the detailed derivations discussed in this section to Appendix.

### 3.1    Differential geometric properties of GBW

To derive the various expressions in Table 1, we provide two strategies, one is by a Riemannian submersion from the general linear group and another is by a Riemannian isometry from the BW Riemannian geometry, $\mathcal{M}_{\mathrm{bw}}$. These claims are formalized in Propositions 3 and 4 respectively.

***Perspective from Riemannian submersion:*** A *Riemannian submersion* [35] between two manifolds is a smooth surjective map where its differential restricted to the horizontal space is isometric (formally defined in Appendix). The general linear group $\mathrm{GL}(n)$ is the set of invertible matrices with the group action of matrix multiplication. When endowed with the standard Euclidean inner product $\langle \cdot, \cdot \rangle_2$, the group becomes a Riemannian manifold, denoted as $\mathcal{M}_{\mathrm{gl}}$. The proposition below introduces a Riemannian submersion from $\mathcal{M}_{\mathrm{gl}}$ to $\mathcal{M}_{\mathrm{gbw}}$.

**Proposition 3.** *The map $\pi : \mathcal{M}_{\mathrm{gl}} \to \mathcal{M}_{\mathrm{gbw}}$ defined as $\pi(\mathbf{P}) = \mathbf{M}^{1/2}\mathbf{P}\mathbf{P}^\top\mathbf{M}^{1/2}$ is a Riemannian submersion, for $\mathbf{P} \in \mathrm{GL}(n)$ and $\mathcal{M}_{\mathrm{gbw}}$ parameterized by $\mathbf{M} \in \mathbb{S}^n_{++}$ as in (3).*

***Perspective from Riemannian isometry:*** A *Riemannian isometry* between two manifolds is a diffeomorphism (i.e., bijective, differentiable, and its inverse is differentiable) that pulls back the Riemannian metric from one to another [35]. We show in the following proposition that there exists a Riemannian isometry between the GBW and BW geometries.

**Table 1.** Summary of expressions for the proposed generalized Bures-Wasserstein (GBW) Riemannian geometry, which is parameterized by $\mathbf{M} \in \mathbb{S}_{++}^d$.

| | |
|---|---|
| Metric | $g_{\mathrm{gbw}}(\mathbf{U}, \mathbf{V}) = \frac{1}{2}\mathrm{tr}(\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{V})$ |
| Distance | $d_{\mathrm{gbw}}^2(\mathbf{X}, \mathbf{Y}) = \mathrm{tr}(\mathbf{M}^{-1}\mathbf{X}) + \mathrm{tr}(\mathbf{M}^{-1}\mathbf{Y}) - 2\mathrm{tr}(\mathbf{X}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1})^{1/2}$ |
| Geodesic | $\gamma(t) = ((1-t)\mathbf{X}^{1/2} + t\mathbf{Y}^{1/2}\mathbf{O})((1-t)\mathbf{X}^{1/2} + t\mathbf{Y}^{1/2}\mathbf{O})^\top$ with $\mathbf{O}$ the orthogonal polar factor of $\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}^{1/2}$. |
| Exp | $\mathrm{Exp}_{\mathbf{X}}(\mathbf{U}) = \mathbf{X} + \mathbf{U} + \mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{X}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{M}$ |
| Log | $\mathrm{Log}_{\mathbf{X}}(\mathbf{Y}) = \mathbf{M}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{Y})^{1/2} + (\mathbf{Y}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{1/2}\mathbf{M} - 2\mathbf{X}$ |
| Connection | $\nabla_\xi \eta = \mathrm{D}_\xi \eta + \{\mathbf{X}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\mathbf{M} + \mathbf{X}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\mathbf{M}\}_{\mathrm{S}} - \{\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\xi\}_{\mathrm{S}} - \{\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\eta\}_{\mathrm{S}}$, where $\{\mathbf{A}\}_{\mathrm{S}} := \frac{1}{2}(\mathbf{A} + \mathbf{A}^\top)$. |
| Min/Max Curvature | $K_{\min}(\pi(\mathbf{P})) = 0$, and $K_{\max}(\pi(\mathbf{P})) = \frac{3}{\sigma_n^2 + \sigma_{n-1}^2}$, where $\sigma_i$ is the $i$-th largest singular value of $\mathbf{P}$, and $\pi(\mathbf{P}) = \mathbf{M}^{1/2}\mathbf{P}\mathbf{P}^\top\mathbf{M}^{1/2}$. |

**Proposition 4.** *Define a map as $\tau(\mathbf{D}) = \mathbf{M}^{-1/2}\mathbf{D}\mathbf{M}^{-1/2}$, for $\mathbf{D} \in \mathbb{S}^n$. Then, the GBW metric can be written as $g_{\mathrm{gbw},\mathbf{X}}(\mathbf{U}, \mathbf{V}) = g_{\mathrm{bw},\tau(\mathbf{X})}(\tau(\mathbf{U}), \tau(\mathbf{V}))$, where the subscript $\mathbf{X}, \tau(\mathbf{X})$ indicates the tangent space. Hence, $\tau : \mathcal{M}_{\mathrm{bw}} \to \mathcal{M}_{\mathrm{gbw}}$ is a Riemannian isometry.*

The proofs of the results in Table 1 are in Appendix and derived from the first perspective of Riemannian submersion, taking inspiration from the analysis in [8,41,42]. In Appendix, we also include various additional developments on the GBW geometry, such as geometric interpolation and barycenter, connection to robust Wasserstein distance and metric learning.

### 3.2 Riemannian optimization with the GBW geometry

Learning over SPD matrices usually concerns optimizing an objective function with respect to the parameter, which is constrained to be SPD. Riemannian optimization is an elegant approach that converts the constrained optimization into an unconstrained problem on manifolds [1,10]. Among the metrics for the SPD matrices, the affine-invariant (AI) metric is seemingly the most popular choice for Riemannian optimization due to its efficiency and convergence guarantees. Recently, however, in [21], the BW metric is shown to be a promising alternative for various learning problems. Below, we derive the expressions for Riemannian gradient and Hessian of an objective function for the GBW geometry.

Riemannian gradient (and Hessian) are generalized gradient (and Hessian) on the tangent space of Riemannian manifolds. The expressions allow to implement various Riemannian optimization methods, using toolboxes like Manopt [11], Pymanopt [53], ROPTLIB [29], etc.

**Proposition 5.** *The Riemannian gradient and Hessian on $\mathcal{M}_{\mathrm{gbw}}$ is derived as* $\mathrm{grad}f(\mathbf{X}) = 2\mathbf{X}\nabla f(\mathbf{X})\mathbf{M} + 2\mathbf{M}\nabla f(\mathbf{X})\mathbf{X}$ *and* $\mathrm{Hess}f(\mathbf{X})[\mathbf{U}] = 4\{\mathbf{M}\nabla^2 f(\mathbf{X})[\mathbf{U}]\mathbf{X}\}_{\mathrm{S}}$

**Table 2.** Riemannian optimization ingredients for the affine-invariant (AI) and Generalized Bures-Wasserstein (GBW) with $\mathbf{M} = \mathbf{X}$ geometries for log-det optimization.

|      | AI | GBW (with $\mathbf{M} = \mathbf{X}$) |
|------|-----|------|
| Exp  | $\mathrm{Exp}_{\mathbf{X}}(\mathbf{U}) = \mathbf{X}\exp(\mathbf{X}^{-1}\mathbf{U})$ | $\mathrm{Exp}_{\mathbf{X}}(\mathbf{U}) = \mathbf{X} + \mathbf{U} + \frac{1}{4}\mathbf{U}\mathbf{X}^{-1}\mathbf{U}$ |
| Grad | $\mathrm{grad}f(\mathbf{X}) = \mathbf{X}\mathbf{C}\mathbf{X} - \mathbf{X}$ | $\mathrm{grad}f(\mathbf{X}) = 4\mathbf{X}\mathbf{C}\mathbf{X} - 4\mathbf{X}$ |
| Hess | $\mathrm{Hess}f(\mathbf{X})[\mathbf{U}] = 2\mathbf{U} + \{\mathbf{U}\mathbf{C}\mathbf{X}\}_{\mathrm{S}}$ | $\mathrm{Hess}f(\mathbf{X})[\mathbf{U}] = 2\mathbf{U} + 2\{\mathbf{U}\mathbf{C}\mathbf{X}\}_{\mathrm{S}}$ |

$$+ 2\{\mathbf{M}\nabla f(\mathbf{X})\mathbf{U}\}_{\mathrm{S}} + 4\{\mathbf{X}\{\nabla f(\mathbf{X})\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\}_{\mathrm{S}}\mathbf{M}\}_{\mathrm{S}} - \{\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathrm{grad}f(\mathbf{X})\}_{\mathrm{S}},$$
where $\nabla f(\mathbf{X}), \nabla^2 f(\mathbf{X})$ represent the Euclidean gradient and Hessian, respectively.

In Appendix, we discuss *geodesic convexity* of functions on the SPD manifold endowed with the GBW metric. It generalizes the discussion in [21].

## 4  Experiments

In this section, we perform experiments showing the benefit of the GBW geometry. The algorithms are implemented in Matlab using the Manopt toolbox [11]. The codes are available on `https://github.com/andyjm3/GBW`.

### 4.1  Log-determinant Riemannian optimization

***Problem formulation:*** Log-determinant (log-det) optimization is common in statistical machine learning, such as for estimating the covariance, with an objective concerning $\min_{\mathbf{X}\in\mathbb{S}^n_{++}} f(\mathbf{X}) = -\log\det(\mathbf{X})$. From [21], optimization with the BW geometry is less well-conditioned compared to the AI geometry. This is because the Riemannian Hessians at optimality are $\mathrm{Hess}_{\mathrm{ai}}f(\mathbf{X}^*)[\mathbf{U}] = \mathbf{U}$ for the AI geometry and $\mathrm{Hess}_{\mathrm{bw}}f(\mathbf{X}^*)[\mathbf{U}] = 4\{(\mathbf{X}^*)^{-1}\mathbf{U}\}_{\mathrm{S}}$ for the BW geometry. This suggests, under the BW geometry, the condition number of Hessian at optimality depends on the solution $\mathbf{X}^*$, while no dependence on $\mathbf{X}^*$ under the AI geometry. Thus, this leads to a poor performance on BW geometry [21].

Here, we show how the GBW geometry helps to address this issue. Specifically, with the GBW geometry, we see from Proposition 5 that by choosing $\mathbf{M} = \mathbf{X}^*$, the Riemannian Hessian is $\mathrm{Hess}_{\mathrm{gbw}}f(\mathbf{X}^*)[\mathbf{U}] = \mathbf{U}$, which becomes well-conditioned (around the optimal solution). This provides the motivation for a choice of $\mathbf{M}$. As the optimal solution $\mathbf{X}^*$ is unknown in optimization problems, choice of $\mathbf{M}$ is not trivial. In practice, one may choose $\mathbf{M} = \mathbf{X}$ dynamically at every or after a few iterations. This strategy corresponds to modifying the GBW geometry dynamically with iterations.

As an example, we consider the following inverse covariance estimation problem [18,27] as $\min_{\mathbf{X}\in\mathbb{S}^n_{++}} f(\mathbf{X}) = -\log\det(\mathbf{X}) + \mathrm{tr}(\mathbf{C}\mathbf{X})$, where $\mathbf{C} \in \mathbb{S}^n_{++}$ is a given SPD matrix. The Euclidean gradient $\nabla f(\mathbf{X}) = -\mathbf{X}^{-1} + \mathbf{C}$ and the Euclidean Hessian $\nabla^2 f(\mathbf{X})[\mathbf{U}] = \mathbf{X}^{-1}\mathbf{U}\mathbf{X}^{-1}$. From the analysis in Appendix, this
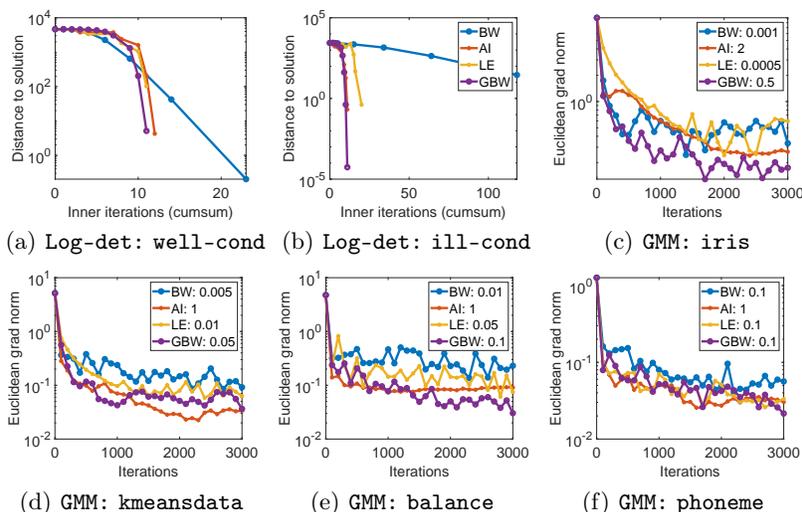
(a) `Log-det: well-cond`    (b) `Log-det: ill-cond`    (c) `GMM: iris`

(d) `GMM: kmeansdata`    (e) `GMM: balance`    (f) `GMM: phoneme`

**Fig. 1.** Figures (a) & (b): Convergence for log-det optimization problem via Riemannian trust region algorithm. Figures (c)-(f): Gaussian mixture model via Riemannian stochastic gradient descent algorithm with optimal initial stepsize. In both the settings, the GBW algorithm outperforms the BW algorithm and performs similar to the AI algorithm. This can be attributed to the choice of $\mathbf{M}$, which offers additional flexibility to the GBW modeling.

problem is geodesic convex and the optimal solution is $\mathbf{X}^* = \mathbf{C}^{-1}$, which we seek to estimate as a direct computation is challenging for ill-conditioned $\mathbf{C}$.

Choosing $\mathbf{M} = \mathbf{X}$ and following derivations in Section 3.2, the expressions for the exponential map, Riemannian gradient, and Hessian under the GBW geometry are shown in Table 2, where we also draw comparisons to the AI geometry. We see that the choice of $\mathbf{M} = \mathbf{X}$ allows GBW to locally approximate the AI geometry up to some constants. For example, the AI exponential map $\mathbf{X}\exp(\mathbf{X}^{-1}\mathbf{U})$ can by approximated by second-order terms as $\mathbf{X} + \mathbf{U} + \frac{1}{2}\mathbf{U}\mathbf{X}^{-1}\mathbf{U}$. This matches the GBW expression up to an additional term $\frac{1}{4}\mathbf{U}\mathbf{X}^{-1}\mathbf{U}$. Overall, the similarity of optimization ingredients help GBW (with $\mathbf{M} = \mathbf{X}$) perform as similar as the AI geometry, which helps to resolve the poor performance of BW for log-det optimization problems observed in [21].

***Experimental setup and results:*** We follow the same settings as in [21] to create problem instances and consider two instances where the condition number of $\mathbf{X}^*$ is 10 (well-conditioned) and 1000 (ill-conditioned). $\mathbf{C}$ is then obtained as $(\mathbf{X}^*)^{-1}$. To compare the convergence performance of optimization methods under the AI, LE, BW, and GBW (with $\mathbf{M} = \mathbf{X}$) geometries, we implement the Riemannian trust region (a second-order solver) with the considered geometries [1,10]. To measure convergence, we use the distance to (theoretical) optimal solution, i.e., $\|\mathbf{X}_t - \mathbf{X}^*\|_2$. We plot this distance against the cumulative inner

iterations that the trust region method takes to solve a particular trust region sub-problem at every iteration. The inner iterations are a good measure to show convergence of trust region algorithms [1, Chapter 7].

From Figures 1(a) & 1(b), we observe the faster convergence with the GBW geometry compared to other geometries regardless of the condition number. In contrast, the BW geometry performs poorly in log-determinant optimization problems as shown in [21]. The GBW geometry effectively resolves the convergence issues with the BW geometry for such settings. Based on our discussion earlier, we see that GBW with $\mathbf{M} = \mathbf{X}$ performs similar to the AI geometry. Empirically, it shows that the GBW geometry effectively bridges the gap between BW and AI geometries for optimization problems.

### 4.2    Gaussian mixture model (GMM)

***Problem formulation:*** We now consider Gaussian density estimation and mixture model problem. Let $\mathbf{x}_i \in \mathbb{R}^d, i = 1, ..., N$, be the given i.i.d. samples. Following [26], we consider a reformulated GMM problem on augmented samples $\mathbf{y}_i^\top = [\mathbf{x}_i^\top; 1] \in \mathbb{R}^{d+1}$. The density of a GMM is parameterized by the augmented covariance matrix $\mathbf{\Sigma} \in \mathbb{R}^{d+1}$. It should be noted that the log-likelihood of Gaussian is geodesic convex under the AI geometry [26] but not under the GBW geometry. However, if we define $\mathbf{S} = \mathbf{\Sigma}^{-1}$ [21], the reparameterized log-likelihood $p_{\mathcal{N}}(\mathbf{Y}; \mathbf{S}) = \sum_{i=1}^{N} \log\left((2\pi)^{1-d/2} \exp(1/2) \det(\mathbf{S})^{1/2} \exp(-\frac{1}{2}\mathbf{y}_i^\top \mathbf{S} \mathbf{y}_i)\right)$ is geodesic convex on $\mathcal{M}_{\mathrm{gbw}}$. Similar trick was employed in [21] to obtained geodesic convex log-likelihood objective for GMM under the BW geometry. Overall, we solve the GMM problem similar as discussed in [26,21].

***Experimental setup and results:*** We consider datasets: `iris`, `kmeansdata`, `balance`, and `phoneme` from Matlab database and Keel database [15]. For comparisons, we implement the Riemannian stochastic gradient descent method [9] as it is widely used in GMM problems [26]. The batch size is set to 50 and we use a decaying stepsize for all the geometries [21]. As discussed in Section 4.1, we set $\mathbf{M} = \mathbf{X}$ at every iteration for optimizing under the GBW geometry. Without access to the optimal solution, the convergence is measured in terms of the Euclidean gradient norm $\|\mathbf{\Sigma}_t \nabla L(\mathbf{\Sigma}_t)\|_2$ for comparability across geometries.

Figures 1(c)-1(f) show convergence along with the best selected initial stepsize. We observe that convergence under the GBW geometry is competitive and clearly outperforms the BW geometry based algorithm.

*Remark 1.* For all the experiments in this section, we simply set $\mathbf{M} = \mathbf{X}$. In general, $\mathbf{M}$ can be learned according to the applications. We demonstrate several examples in Appendix.

## 5    Conclusion

In this paper, we propose a Riemannian geometry that generalizes the recently introduced Bures-Wasserstein geometry for SPD matrices. This generalized geometry has natural connections to the orthogonal Procrustes problem as well

as to the optimal transport theory, and still possesses the properties of the Bures-Wasserstein geometry (which is a special case). The new geometry is shown to be parameterized by a SPD matrix $\mathbf{M}$. This offers necessary flexibility in applications. Experiments show that learning of $\mathbf{M}$ leads to better modeling in applications.

## References

1. P.-A. Absil, R. Mahony, and R. Sepulchre. *Optimization algorithms on matrix manifolds.* Princeton University Press, 2008.
2. Martial Agueh and Guillaume Carlier. Barycenters in the Wasserstein space. *SIAM Journal on Mathematical Analysis*, 43(2):904–924, 2011.
3. Pedro C Álvarez-Esteban, E Del Barrio, JA Cuesta-Albertos, and C Matrán. A fixed-point approach to barycenters in Wasserstein space. *Journal of Mathematical Analysis and Applications*, 441(2):744–762, 2016.
4. Vincent Arsigny, Pierre Fillard, Xavier Pennec, and Nicholas Ayache. Log-Euclidean metrics for fast and simple calculus on diffusion tensors. *Magnetic Resonance in Medicine: An Official Journal of the International Society for Magnetic Resonance in Medicine*, 56(2):411–421, 2006.
5. Richard Bellman. Some inequalities for the square root of a positive definite matrix. *Linear Algebra and its applications*, 1(3):321–324, 1968.
6. Arthur L Besse. *Einstein manifolds.* Springer Science & Business Media, 2007.
7. Rajendra Bhatia. *Positive definite matrices.* Princeton university press, 2009.
8. Rajendra Bhatia, Tanvi Jain, and Yongdo Lim. On the Bures-Wasserstein distance between positive definite matrices. *Expositiones Mathematicae*, 37(2):165–191, 2019.
9. Silvere Bonnabel. Stochastic gradient descent on riemannian manifolds. *IEEE Transactions on Automatic Control*, 58(9):2217–2229, 2013.
10. N. Boumal. *An introduction to optimization on smooth manifolds.* Available online, Aug, 2020.
11. N. Boumal, B. Mishra, P.-A. Absil, and R. Sepulchre. Manopt, a Matlab toolbox for optimization on manifolds. *Journal of Machine Learning Research*, 15(1):1455–1459, 2014.
12. Daniel A Brooks, Olivier Schwander, Frédéric Barbaresco, Jean-Yves Schneider, and Matthieu Cord. Exploring complex time-series representations for Riemannian machine learning of radar data. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2019.
13. Anoop Cherian and Suvrit Sra. Riemannian dictionary learning and sparse coding for positive definite matrices. *IEEE transactions on neural networks and learning systems*, 28(12):2859–2871, 2016.
14. Pinar Demetci, Rebecca Santorella, Bjorn Sandstede, William Stafford Noble, and Ritambhara Singh. Gromov-wasserstein optimal transport to align single-cell multi-omics data. *BioRxiv*, 2020.
15. J Derrac, S Garcia, L Sanchez, and F Herrera. Keel data-mining software tool: Data set repository, integration of algorithms and experimental analysis framework. *J. Mult. Valued Logic Soft Comput*, 17, 2015.
16. R. M. Dudley. The speed of mean Glivenko-Cantelli convergence. *The Annals of Mathematical Statistics*, 40(1):40–50, 1969.
17. Nicolas Fournier and Arnaud Guillin. On the rate of convergence in Wasserstein distance of the empirical measure. *Probability Theory and Related Fields*, 162(3):707–738, 2015.

18. Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2008.
19. Alexandre Gramfort, Gabriel Peyré, and Marco Cuturi. Fast optimal transport averaging of neuroimaging data. In *International Conference on Information Processing in Medical Imaging*, 2015.
20. Matthieu Guillaumin, Jakob Verbeek, and Cordelia Schmid. Is that you? metric learning approaches for face identification. In *International Conference on Computer Vision*, 2009.
21. Andi Han, Bamdev Mishra, Pratik Jawanpuria, and Junbin Gao. On Riemannian optimization over positive definite matrices with the Bures-Wasserstein geometry. In *Advances in Neural Information Processing Systems*, 2021.
22. Mehrtash Harandi, Mathieu Salzmann, and Richard Hartley. Joint dimensionality reduction and metric learning: A geometric take. In *International Conference on Machine Learning*, 2017.
23. Mehrtash T Harandi, Richard Hartley, Brian Lovell, and Conrad Sanderson. Sparse coding on symmetric positive definite manifolds using bregman divergences. *IEEE transactions on neural networks and learning systems*, 27(6):1294–1306, 2015.
24. Mehrtash T Harandi, Mathieu Salzmann, and Richard Hartley. From manifold to manifold: Geometry-aware dimensionality reduction for SPD matrices. In *European Conference on Computer Vision*, 2014.
25. Inbal Horev, Florian Yger, and Masashi Sugiyama. Geometry-aware principal component analysis for symmetric positive definite matrices. In *Asian Conference on Machine Learning*, 2016.
26. Reshad Hosseini and Suvrit Sra. An alternative to EM for Gaussian mixture models: batch and stochastic Riemannian optimization. *Mathematical Programming*, 181(1):187–223, 2020.
27. Cho-Jui Hsieh, Inderjit Dhillon, Pradeep Ravikumar, and Mátyás Sustik. Sparse inverse covariance matrix estimation using quadratic approximation. In *Advances in neural information processing systems*, 2011.
28. Minhui Huang, Shiqian Ma, and Lifeng Lai. Projection robust Wasserstein barycenters. In *International Conference on Machine Learning*, 2021.
29. Wen Huang, P-A Absil, Kyle A Gallivan, and Paul Hand. Roptlib: an object-oriented c++ library for optimization on riemannian manifolds. *ACM Transactions on Mathematical Software (TOMS)*, 44(4):1–21, 2018.
30. Zhiwu Huang, Ruiping Wang, Xianqiu Li, Wenxian Liu, Shiguang Shan, Luc Van Gool, and Xilin Chen. Geometry-aware similarity learning on SPD manifolds for visual recognition. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(10):2513–2523, 2017.
31. Zhiwu Huang, Ruiping Wang, Shiguang Shan, Xianqiu Li, and Xilin Chen. Log-Euclidean metric learning on symmetric positive definite manifold with application to image set classification. In *International conference on Machine Learning*, 2015.
32. Minyoung Kim, Sanjiv Kumar, Vladimir Pavlovic, and Henry Rowley. Face tracking and recognition with visual constraints in real-world videos. In *Conference on Computer Vision and Pattern Recognition*, 2008.
33. Serge Lang. *Differential and Riemannian manifolds*. Springer Science & Business Media, 2012.
34. Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
35. John M Lee. *Riemannian manifolds: an introduction to curvature*, volume 176. Springer Science & Business Media, 2006.

36. John M Lee. *Introduction to Riemannian manifolds*. Springer, 2018.
37. Bastian Leibe and Bernt Schiele. Analyzing appearance and contour based methods for object categorization. In *Conference on Computer Vision and Pattern Recognition*, 2003.
38. Zhenhua Lin. Riemannian geometry of symmetric positive definite matrices via cholesky decomposition. *SIAM Journal on Matrix Analysis and Applications*, 40(4):1353–1370, 2019.
39. Fabien Lotte, Laurent Bougrain, Andrzej Cichocki, Maureen Clerc, Marco Congedo, Alain Rakotomamonjy, and Florian Yger. A review of classification algorithms for eeg-based brain–computer interfaces: a 10 year update. *Journal of neural engineering*, 15(3):031005, 2018.
40. Sridhar Mahadevan, Bamdev Mishra, and Shalini Ghosh. A unified framework for domain adaptation using metric learning on manifolds. In *European Conference on Machine Learning and Knowledge Discovery in Databases*, 2019.
41. Luigi Malagò, Luigi Montrucchio, and Giovanni Pistone. Wasserstein Riemannian geometry of Gaussian densities. *Information Geometry*, 1(2):137–179, 2018.
42. Estelle Massart, Julien M Hendrickx, and P-A Absil. Curvature of the manifold of fixed-rank positive-semidefinite matrices endowed with the Bures-Wasserstein metric. In *International Conference on Geometric Science of Information*, 2019.
43. Bamdev Mishra and Rodolphe Sepulchre. Riemannian preconditioning. *SIAM Journal on Optimization*, 26(1):635–660, 2016.
44. Barrett O'Neill. The fundamental equations of a submersion. *Michigan Mathematical Journal*, 13(4):459–469, 1966.
45. François-Pierre Paty and Marco Cuturi. Subspace robust Wasserstein distances. In *International Conference on Machine Learning*, 2019.
46. Xavier Pennec, Pierre Fillard, and Nicholas Ayache. A Riemannian framework for tensor computing. *International Journal of computer vision*, 66(1):41–66, 2006.
47. G. Peyré and M. Cuturi. Computational optimal transport. *Foundations and Trends in Machine Learning*, 11(5-6):355–607, 2019.
48. Julien Rabin, Gabriel Peyré, Julie Delon, and Marc Bernot. Wasserstein barycenter and its application to texture mixing. In *International Conference on Scale Space and Variational Methods in Computer Vision*, 2011.
49. Boris Shustin and Haim Avron. Preconditioned Riemannian optimization on the generalized Stiefel manifold. *arXiv:1902.01635*, 2019.
50. Suvrit Sra. A new metric on the manifold of kernel matrices with application to matrix geometric means. *Advances in Neural Information Processing Systems*, 2012.
51. Suvrit Sra. Positive definite matrices and the S-divergence. *Proceedings of the American Mathematical Society*, 144(7):2787–2797, 2016.
52. Suvrit Sra and Reshad Hosseini. Conic geometric optimization on the manifold of positive definite matrices. *SIAM Journal on Optimization*, 25(1):713–739, 2015.
53. J. Townsend, N. Koep, and S. Weichwald. Pymanopt: A python toolbox for optimization on manifolds using automatic differentiation. *Journal of Machine Learning Research*, 17(137):1–5, 2016.
54. Jesse van Oostrum. Bures-Wasserstein geometry. *arXiv:2001.08056*, 2020.
55. J. Weed and F. Bach. Sharp asymptotic and finite-sample rates of convergence of empirical measures in Wasserstein distance. *Bernoulli*, 25(4A):2620–2648, 2019.
56. Pourya Zadeh, Reshad Hosseini, and Suvrit Sra. Geometric mean metric learning. In *International Conference on Machine Learning*, 2016.
57. Hongyi Zhang and Suvrit Sra. First-order methods for geodesically convex optimization. In *Conference on Learning Theory*, 2016.

## A  Additional results and proofs for Section 2

### A.1  Proof of Proposition 1

*Proof (Proof of Proposition 1).* First we have

$$
\min_{\mathbf{O}\in O(n)} \|\mathbf{X}^{1/2} - \mathbf{Y}^{1/2}\mathbf{O}\|_{\mathbf{M}^{-1}}^2
$$
$$
= \operatorname{tr}(\mathbf{M}^{-1}\mathbf{X}) + \operatorname{tr}(\mathbf{M}^{-1}\mathbf{Y}) - 2 \max_{\mathbf{O}\in O(n)} \operatorname{tr}(\mathbf{M}^{-1}\mathbf{X}^{1/2}\mathbf{O}\mathbf{Y}^{1/2}). \qquad (5)
$$

And the minimum of (5) is attained when $\mathbf{O}$ is the orthogonal polar factor of $\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}^{1/2}$, which is $\mathbf{O} = \mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}^{1/2}(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{-1/2}$, as proved in [8]. Substituting the expression of $\mathbf{O}$ in (5) completes the proof.

### A.2  Proof of Proposition 2

Before we proceed to prove Proposition 2, we provide an essential lemma, which generalizes [8, Theorem 2].

**Lemma 1.** *Define* $\tilde{F}(\mathbf{X}, \mathbf{Y}) = \operatorname{tr}(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{1/2}$. *Then for any* $\mathbf{X}, \mathbf{Y} \in \mathbb{S}_{++}^n$,

1) $\tilde{F}(\mathbf{X}, \mathbf{Y}) = \min_{\mathbf{A}\in\mathbb{S}_{++}^n} \frac{1}{2}\operatorname{tr}(\mathbf{X}\mathbf{A} + \mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{A}^{-1})$.
2) $\tilde{F}(\mathbf{X}, \mathbf{Y}) = \min_{\mathbf{A}\in\mathbb{S}_{++}^n} \sqrt{\operatorname{tr}(\mathbf{X}\mathbf{A})\operatorname{tr}(\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{A}^{-1})}$.

*Proof.* Following [8], the proof proceeds by analyzing the first-order stationary conditions, where we replace $\mathbf{Y}$ with $\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}$.

*Proof (Proof of Proposition 2).* We have $\mathbf{X} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$ and $\mathbf{Y} = \mathbb{E}[\mathbf{y}\mathbf{y}^\top]$.

$$
\tilde{W}^2(\mu, \nu) = \inf_{\mathbf{x}\sim\mu, \mathbf{y}\sim\nu} \mathbb{E}[\mathbf{x}^\top\mathbf{M}^{-1}\mathbf{x} + \mathbf{y}^\top\mathbf{M}^{-1}\mathbf{y} - 2\mathbf{x}^\top\mathbf{M}^{-1}\mathbf{y}]
$$
$$
= \inf_{\mathbf{x}\sim\mu, \mathbf{y}\sim\nu} \operatorname{tr}(\mathbf{M}^{-1}\mathbf{X}) + \operatorname{tr}(\mathbf{M}^{-1}\mathbf{Y}) - 2\operatorname{tr}(\mathbf{M}^{-1}\mathbb{E}[\mathbf{y}\mathbf{x}^\top])]
$$
$$
= \operatorname{tr}(\mathbf{M}^{-1}\mathbf{X}) + \operatorname{tr}(\mathbf{M}^{-1}\mathbf{Y}) - \sup_{\mathbf{K}:\mathbf{\Sigma}\succeq\mathbf{0}} 2\operatorname{tr}(\mathbf{M}^{-1}\mathbf{K}^\top),
$$

where $\mathbf{K}$ is the covariance between $\mathbf{x}, \mathbf{y}$ such that the joint covariance matrix

$$
\mathbf{\Sigma} = \mathbb{E}\begin{bmatrix} \mathbf{x}\mathbf{x}^\top & \mathbf{x}\mathbf{y}^\top \\ \mathbf{x}\mathbf{y}^\top & \mathbf{y}\mathbf{y}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{X} & \mathbf{K} \\ \mathbf{K}^\top & \mathbf{Y} \end{bmatrix} \succeq \mathbf{0}.
$$

Two necessary and sufficient conditions for $\mathbf{\Sigma} \succeq \mathbf{0}$ are (i) $\mathbf{X} \succeq \mathbf{K}\mathbf{Y}^{-1}\mathbf{K}^\top$ and (ii) $\mathbf{K} = \mathbf{X}^{1/2}\mathbf{C}\mathbf{Y}^{1/2}$ for some contraction $\mathbf{C}$, i.e., $\|\mathbf{C}\|_2 \leq 1$ [7]. Hence, $\operatorname{tr}(\mathbf{K}) \leq \|\mathbf{X}^{1/2}\|_2\|\mathbf{Y}^{1/2}\|_2 = \sqrt{\operatorname{tr}(\mathbf{X})\operatorname{tr}(\mathbf{Y})}$. Also, for any $\mathbf{A} \in \mathbb{S}_{++}^n$, the block diagonal matrix

$$
\mathbf{P} = \begin{bmatrix} \mathbf{A}^{1/2} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}^{-1/2}\mathbf{M}^{-1} \end{bmatrix} \in \mathrm{M}(2n)
$$

Then

$$\mathbf{P}\begin{bmatrix} \mathbf{X} & \mathbf{K} \\ \mathbf{K}^\top & \mathbf{Y} \end{bmatrix}\mathbf{P}^\top = \begin{bmatrix} \mathbf{A}^{1/2}\mathbf{X}\mathbf{A}^{1/2} & \mathbf{A}^{1/2}\mathbf{K}\mathbf{M}^{-1}\mathbf{A}^{-1/2} \\ \mathbf{A}^{-1/2}\mathbf{M}^{-1}\mathbf{K}^\top\mathbf{A}^{1/2} & \mathbf{A}^{-1/2}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{A}^{-1/2} \end{bmatrix} \succeq \mathbf{0}.$$

This leads to

$$\begin{aligned}
\operatorname{tr}(\mathbf{M}^{-1}\mathbf{K}^\top) &= \operatorname{tr}(\mathbf{A}^{-1/2}\mathbf{M}^{-1}\mathbf{K}^\top\mathbf{A}^{1/2}) \\
&\leq \sqrt{\operatorname{tr}(\mathbf{A}^{1/2}\mathbf{X}\mathbf{A}^{1/2})\operatorname{tr}(\mathbf{A}^{-1/2}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{A}^{-1/2})} \\
&= \sqrt{\operatorname{tr}(\mathbf{X}\mathbf{A})\operatorname{tr}(\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{A}^{-1})}.
\end{aligned}$$

Hence choosing $\mathbf{K} = (\mathbf{X}\mathbf{Y})^{1/2}$, we can show $\mathbf{K}\mathbf{Y}^{-1}\mathbf{K}^\top = \mathbf{X}$ [8] and

$$\max_{\mathbf{K}:\boldsymbol{\Sigma}\succeq\mathbf{0}} \operatorname{tr}(\mathbf{M}^{-1}\mathbf{K}^\top) = \tilde{F}(\mathbf{X},\mathbf{Y}) = \min_{\mathbf{A}\in\mathbb{S}_{++}^n} \sqrt{\operatorname{tr}(\mathbf{X}\mathbf{A})\operatorname{tr}(\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{A}^{-1})}.$$

It thus follows that $\tilde{W}^2(\mu,\nu) = d_{\mathrm{gbw}}^2(\mathbf{X},\mathbf{Y})$.

### A.3  Additional results: optimal transport plan for GBW distance

Next, we derive the optimal transport plan corresponding to the generalized Wasserstein distance as follows, where we use the notation $\mathbf{A}\#\mathbf{B}$ to represent the matrix geometric mean under the affine-invariant metric, i.e., $\mathbf{A}\#\mathbf{B} = \mathbf{A}^{1/2}(\mathbf{A}^{-1/2}\mathbf{B}\mathbf{A}^{-1/2})^{1/2}\mathbf{A}^{1/2} = \mathbf{A}(\mathbf{A}^{-1}\mathbf{B})^{1/2} = (\mathbf{A}\mathbf{B}^{-1})^{1/2}\mathbf{B}$.

**Proposition 6.** *Let* $\mathbf{x},\mathbf{y} \in \mathbb{R}^n$ *be random Gaussian vectors with zero mean and covariance matrices* $\mathbf{X},\mathbf{Y} \in \mathbb{S}_{++}^n$ *respectively. The optimal transport plan from* $\mathbf{x}$ *to* $\mathbf{y}$ *under the Mahalanobis distance is* $\mathbf{T} = \mathbf{M}(\mathbf{X}^{-1}\#(\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}))$.

*Proof (Proof of Proposition 6).* For any $\mathbf{P} \in \mathrm{M}(n)$ as a transport plan,

$$\mathbb{E}\|\mathbf{M}^{-1/2}\mathbf{x} - \mathbf{P}\mathbf{M}^{-1/2}\mathbf{x}\|_2^2$$
$$= \operatorname{tr}(\mathbf{M}^{-1}\mathbf{X}) + \operatorname{tr}(\mathbf{P}\mathbf{M}^{-1/2}\mathbf{X}\mathbf{M}^{-1/2}\mathbf{P}^\top) - 2\operatorname{tr}(\mathbf{X}^{1/2}\mathbf{M}^{-1/2}\mathbf{P}\mathbf{M}^{-1/2}\mathbf{X}^{1/2}). \quad (6)$$

By comparing (6) to $d_{\mathrm{gbw}}^2(\mathbf{X},\mathbf{Y})$, we set

$$\mathbf{X}^{1/2}\mathbf{M}^{-1/2}\mathbf{P}\mathbf{M}^{-1/2}\mathbf{X}^{1/2} = (\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{1/2}$$

which gives an expression of $\mathbf{P}$ as

$$\begin{aligned}
\mathbf{P} &= \mathbf{M}^{1/2}\mathbf{X}^{-1/2}(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{1/2}\mathbf{X}^{-1/2}\mathbf{M}^{1/2} \\
&= \mathbf{M}^{1/2}\big(\mathbf{X}^{-1}\#(\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1})\big)\mathbf{M}^{1/2}. \quad (7)
\end{aligned}$$

From this result, we have

$$\begin{aligned}
&\operatorname{tr}(\mathbf{P}\mathbf{M}^{-1/2}\mathbf{X}\mathbf{M}^{-1/2}\mathbf{P}^\top) \\
&= \operatorname{tr}\Big(\mathbf{M}^{1/2}\big(\mathbf{X}^{-1}\#(\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1})\big)\mathbf{X}\big(\mathbf{X}^{-1}\#(\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1})\big)\mathbf{M}^{1/2}\Big) \\
&= \operatorname{tr}\Big(\mathbf{M}^{1/2}\mathbf{X}^{-1}(\mathbf{X}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1})^{1/2}\mathbf{X}\mathbf{X}^{-1}(\mathbf{X}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1})^{1/2}\mathbf{M}^{1/2}\Big) \\
&= \operatorname{tr}(\mathbf{M}^{-1}\mathbf{Y}),
\end{aligned}$$

where we use the property of matrix geometric mean. This suggests the definition of $\mathbf{P}$ is the optimal transport map under the Euclidean distance. Combining (7) with (6) shows

$$\mathbb{E}\|\mathbf{M}^{-1/2}\mathbf{x} - \mathbf{M}^{1/2}\big(\mathbf{X}^{-1}\#(\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1})\big)\mathbf{x}\|_2^2$$
$$= \mathbb{E}\|\mathbf{x} - \mathbf{M}\big(\mathbf{X}^{-1}\#(\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1})\big)\mathbf{x}\|_{\mathbf{M}^{-1}}^2$$
$$= d_{\mathrm{gbw}}^2(\mathbf{X}, \mathbf{Y}).$$

We can, thus, define the transport plan as $\mathbf{T_{X\to Y}} := \mathbf{M}\big(\mathbf{X}^{-1}\#(\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1})\big)$ and denote $\mathbf{y} = \mathbf{T_{X\to Y}}\mathbf{x}$, which is a Gaussian random vector with covariance

$$\mathbf{T_{X\to Y}}\mathbf{X}\mathbf{T_{X\to Y}^\top} = \mathbf{M}\big(\mathbf{X}^{-1}\#(\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1})\big)\mathbf{X}\big(\mathbf{X}^{-1}\#(\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1})\big)\mathbf{M} = \mathbf{Y}.$$

Thus, $\mathbb{E}\|\mathbf{x} - \mathbf{y}\|_{\mathbf{M}^{-1}}^2 = d_{\mathrm{gbw}}^2(\mathbf{X}, \mathbf{Y})$ where $\mathbf{y} = \mathbf{T_{X\to Y}}\mathbf{x}$. From Proposition 2, we see $\mathbf{T_{X\to Y}}$ is the optimal transport plan from $\mathbf{X}$ to $\mathbf{Y}$ under the Mahalanobis distance.

# B    Additional results and proofs for Section 3.1

## B.1    Riemannian distance, geodesics, exponential map, and logarithm map

**Proposition 7.** *The Riemannian distance on $\mathcal{M}_{\mathrm{gbw}}$ is derived as $d_{\mathrm{gbw}}(\mathbf{X}, \mathbf{Y}) = (\mathrm{tr}(\mathbf{M}^{-1}\mathbf{X}) + \mathrm{tr}(\mathbf{M}^{-1}\mathbf{Y}) - 2\mathrm{tr}(\mathbf{X}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1})^{1/2})^{1/2}$.*

**Proposition 8.** *A geodesic on $\mathcal{M}_{\mathrm{gbw}}$ between any $\mathbf{X}, \mathbf{Y} \in \mathbb{S}_{++}^n$ is given by $\gamma(t) = (\pi \circ c)(t)$, where $c(t) = (1-t)\mathbf{M}^{-1/2}\mathbf{X}^{1/2} + t\mathbf{M}^{-1/2}\mathbf{Y}^{1/2}\mathbf{O}$. Here, $\mathbf{O}$ is the orthogonal polar factor of $\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}^{1/2}$.*

The geodesic in Proposition 8 can be simplified as

$$\gamma(t) = \psi(t)\psi(t)^\top = \big((1-t)\mathbf{X}^{1/2} + t\mathbf{Y}^{1/2}\mathbf{O}\big)\big((1-t)\mathbf{X}^{1/2} + t\mathbf{Y}^{1/2}\mathbf{O}\big)^\top, \quad (8)$$

which coincides with the geodesic of the BW geometry except that $\mathbf{O}$ is now the orthogonal polar factor of $\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}^{1/2}$ rather than $\mathbf{Y}^{1/2}\mathbf{X}^{1/2}$ as for BW.

**Proposition 9.** *The Riemannian exponential map associated with the generalized BW metric is $\mathrm{Exp}_{\mathbf{X}}(t\mathbf{U}) = \mathbf{X} + t\mathbf{U} + t^2\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{X}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{M}$. The neighbourhood $\mathcal{X} := \{\mathbf{M} + t\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{M} \in \mathbb{S}_{++}^n\}$ is a totally normal neighbourhood where exponential map is a diffeomorphism with logarithm map $\mathrm{Log}_{\mathbf{X}}(\mathbf{Y}) = \mathbf{M}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{Y})^{1/2} + (\mathbf{Y}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{1/2}\mathbf{M} - 2\mathbf{X}$.*

We remark that the exponential map is only invertible in the neighbourhood $\mathcal{X}$, where $t$ is chosen sufficiently small. This makes $\mathcal{M}_{\mathrm{gbw}}$ a geodesic incomplete manifold, similar to $\mathcal{M}_{\mathrm{bw}}$ [41].

## B.2   Levi-Civita connection and sectional curvature

The Levi-Civita connection (Levi-Civita derivative) of a vector field on manifold $\mathcal{M}$ is the unique covariant derivative that satisfies (1) torsion-free property, and (2) metric compatibility (formal definitions are provided as supplementary). Let $\mathfrak{X}(\mathcal{M})$ be the space of vector fields on the Riemannian manifold $(\mathcal{M}, g)$ and denote $\{\mathbf{A}\}_\mathrm{S} := (\mathbf{A} + \mathbf{A})/2$, for $\mathbf{A} \in \mathbb{R}^{n \times n}$.

**Proposition 10.** *The Levi-Civita connection with the GBW geometry is* $\nabla_\xi \eta = \mathrm{D}_\xi \eta + \{\mathbf{X}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\mathbf{M} + \mathbf{X}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\mathbf{M}\}_\mathrm{S} - \{\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\xi\}_\mathrm{S} - \{\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\eta\}_\mathrm{S}$, *for any* $\xi, \eta \in \mathfrak{X}(\mathcal{M})$.

Sectional curvature measures the curvature locally around a point $\mathbf{X}$, which is defined geometrically as the Gaussian curvature of a 2-dimensional subspace of $T_\mathbf{X}\mathcal{M}_\mathrm{gbw}$.

**Proposition 11.** *Let* $U, V \in \mathfrak{X}(\mathcal{M}_\mathrm{gbw})$ *be two (linearly independent) vector fields. Let* $\tilde{U}(\mathbf{P}) = \mathbf{M}^{1/2}\mathcal{L}_{\mathbf{M},\pi(\mathbf{P})}[U(\pi(\mathbf{P}))]\mathbf{M}^{1/2}\mathbf{P}$, *for* $\mathbf{P} \in \mathcal{M}_\mathrm{gl}$ *and similarly for* $\tilde{V}$. *Suppose* $\tilde{U}(\mathbf{P})$, $\tilde{V}(\mathbf{P})$ *are orthonormal on* $T_\mathbf{P}\mathcal{M}_\mathrm{gl}$. *Then, the sectional curvature of the subspace spanned by* $U(\pi(\mathbf{P})), V(\pi(\mathbf{P}))$ *is*

$$K(U(\pi(\mathbf{P})), V(\pi(\mathbf{P}))) = \sum_{i,j} \frac{3\mathbf{C}_{ij}^2}{\sigma_j^2(\sigma_i\sigma_j^{-1} + \sigma_i^{-1}\sigma_j)^2},$$

*where* $\mathbf{C} = \mathbf{V}^\top(\tilde{V}(\mathbf{P})^\top\tilde{U}(\mathbf{P}) - \tilde{U}(\mathbf{P})^\top\tilde{V}(\mathbf{P}))\mathbf{V}$ *and the singular value decomposition gives* $\mathbf{P} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ *with* $\mathbf{\Sigma} = \mathrm{diag}(\sigma_1, \ldots, \sigma_n)$, $\sigma_1 \geq \cdots \geq \sigma_n$.

The sectional curvature of GBW geometry is shown to be non-negative, regardless of the choice of $\mathbf{M}$. The bounds are shown below.

**Proposition 12.** *Under the same settings as Proposition 11, the minimum sectional curvature is zero and the maximum is* $3/(\sigma_n^2 + \sigma_{n-1}^2)$.

Although sharing the minimum curvature with the BW geometry [42], the maximum curvature of GBW geometry is affected by the choice of $\mathbf{M}$ from the definition of Riemannian submersion $\pi$ defined in Proposition 3.

## B.3   Geometric interpolation and barycenter

With the proposed GBW geometry, we are interested to study the properties of interpolation between two or more SPD matrices on the manifold. This has implications in various applications, such as diffusion tensor imaging (DTI) [46,7].

First, we show that the geodesic interpolation between $\mathbf{X}, \mathbf{Y} \in \mathcal{M}_\mathrm{gbw}$ satisfies an operator inequality (shown in supplementary), which is $\gamma(t) \preceq (1-t)\mathbf{X} + t\mathbf{Y}$, where $\gamma(0) = \mathbf{X}$, $\gamma(1) = \mathbf{Y}$ and $\preceq$ denotes the Löwner partial order. One immediate implication is $\log\det(\gamma(t)) \leq \log\det((1-t)\mathbf{X} + t\mathbf{Y})$ that suggests a smaller swelling effect compared to the Euclidean interpolation for DTI applications. See more discussions in the supplementary.

We further study the interpolation for multiple SPD matrices $\{\mathbf{X}_l\}_{l=1}^N$, also known as the barycenter problem on $\mathcal{M}_{\mathrm{gbw}}$, i.e., $\min_{\mathbf{A} \in \mathbb{S}_{++}^n} \sum_{l=1}^N w_l d_{\mathrm{gbw}}^2(\mathbf{X}_l, \mathbf{A})$, where the weights $\sum_l w_l = 1$. We show in the supplementary that there exists a unique solution to the problem and provide a fixed point iteration to compute the barycenter [8].

### B.4   Proof of Proposition 3

First we recall a *smooth submersion* is a smooth map $\pi : (\mathcal{M}, g) \to (\mathcal{N}, h)$ from Riemannian manifold $(\mathcal{M}, g)$ to $(\mathcal{N}, h)$ such that its differential $\mathrm{D}\pi(x) : T_x\mathcal{M} \to T_{\pi(x)}\mathcal{N}$ is surjective for any $x \in \mathcal{M}$. Every tangent space $T_x\mathcal{M}$ can be decomposed as $T_x\mathcal{M} = \mathcal{V}_x \oplus \mathcal{H}_x = \mathrm{Ker}(\mathrm{D}\pi(x)) \oplus \mathrm{Ker}(\mathrm{D}\pi(x))^\perp$, where $\mathrm{Ker}(f)$ denotes the kernel of a map and $\oplus$ is the direct sum. We respectively call $\mathcal{V}_x, \mathcal{H}_x$ as the vertical and horizontal subspaces. The map $\pi$ is called a *Riemannian submersion* if it is a smooth submersion and its differential restricted to the horizontal space, $\mathrm{D}\pi(x) : \mathcal{H}_x \to T_{\pi(x)}\mathcal{N}$ is isometric for any $x \in \mathcal{M}$, i.e. $g_x(u, v) = h_{\pi(x)}(\mathrm{D}\pi(x)[u], \mathrm{D}\pi(x)[v])$.

*Proof (Proof of Proposition 3).* Note that the tangent space of $\mathcal{M}_{\mathrm{gl}}$, $T_{\mathbf{P}}\mathcal{M}_{\mathrm{gl}}$, is the space of $\mathbb{R}^{n \times n}$. The differential of $\pi(\mathbf{P})$ in the direction $\mathbf{U} \in \mathbb{R}^{n \times n}$ is given by $\mathrm{D}\pi(\mathbf{P})[\mathbf{U}] = \mathbf{M}^{1/2}\mathbf{U}\mathbf{P}^\top\mathbf{M}^{1/2} + \mathbf{M}^{1/2}\mathbf{P}\mathbf{U}^\top\mathbf{M}^{1/2}$. We then derive the kernel of $\mathrm{D}\pi(\mathbf{P})$ (vertical space $\mathcal{V}_{\mathbf{P}}$) and the orthogonal complement of the kernel (horizontal space $\mathcal{H}_{\mathbf{P}}$) as

$$
\begin{aligned}
\mathrm{Ker}(\mathrm{D}\pi(\mathbf{P})) &= \{\mathbf{U} : \mathrm{D}\pi(\mathbf{P})[\mathbf{U}] = \mathbf{0}\} \\
&= \{\mathbf{U} = \mathbf{M}^{-1/2}\mathbf{K}\mathbf{M}^{-1/2}\mathbf{P}^{-\top} : \mathbf{K} \text{ is skew-symmetric}\}, \quad (9) \\
\mathrm{Ker}(\mathrm{D}\pi(\mathbf{P}))^\perp &= \{\mathbf{V} : \mathrm{tr}(\mathbf{V}^\top\mathbf{M}^{-1/2}\mathbf{K}\mathbf{M}^{-1/2}\mathbf{P}^{-\top}) = \mathbf{0}\} \\
&= \{\mathbf{V} = \mathbf{M}^{1/2}\mathbf{S}\mathbf{M}^{1/2}\mathbf{P} : \mathbf{S} \in \mathbb{S}^n\}. \quad (10)
\end{aligned}
$$

It is clear that $\pi$ is a smooth submersion. Now, we only need to verify that it also satisfies the isometry property. For any $\mathbf{S}, \mathbf{H} \in \mathbb{S}^n$, $\mathbf{M}^{1/2}\mathbf{S}\mathbf{M}^{1/2}\mathbf{P}, \mathbf{M}^{1/2}\mathbf{H}\mathbf{M}^{1/2}\mathbf{P} \in \mathcal{H}_{\mathbf{P}}$, and

$$
\begin{aligned}
\mathrm{D}\pi(\mathbf{P})[\mathbf{M}^{1/2}\mathbf{S}\mathbf{M}^{1/2}\mathbf{P}] &= \mathbf{M}\mathbf{S}\mathbf{M}^{1/2}\mathbf{P}\mathbf{P}^\top\mathbf{M}^{1/2} + \mathbf{M}^{1/2}\mathbf{P}\mathbf{P}^\top\mathbf{M}^{1/2}\mathbf{S}\mathbf{M} \\
&= \mathbf{M}\mathbf{S}\pi(\mathbf{P}) + \pi(\mathbf{P})\mathbf{S}\mathbf{M} \\
\mathrm{D}\pi(\mathbf{P})[\mathbf{M}^{1/2}\mathbf{H}\mathbf{M}^{1/2}\mathbf{P}] &= \mathbf{M}\mathbf{H}\mathbf{M}^{1/2}\mathbf{P}\mathbf{P}^\top\mathbf{M}^{1/2} + \mathbf{M}^{1/2}\mathbf{P}\mathbf{P}^\top\mathbf{M}^{1/2}\mathbf{H}\mathbf{M} \\
&= \mathbf{M}\mathbf{H}\pi(\mathbf{P}) + \pi(\mathbf{P})\mathbf{H}\mathbf{M}.
\end{aligned}
$$

The inner product at $\pi(\mathbf{P})$ is given by

$$
\begin{aligned}
&\langle \mathrm{D}\pi(\mathbf{P})[\mathbf{M}^{1/2}\mathbf{S}\mathbf{M}^{1/2}\mathbf{P}], \mathrm{D}\pi(\mathbf{P})[\mathbf{M}^{1/2}\mathbf{H}\mathbf{M}^{1/2}\mathbf{P}] \rangle_{\mathrm{gbw}} \\
=&\frac{1}{2}\mathrm{tr}(\mathcal{L}_{\pi(\mathbf{P}),\mathbf{M}}[\mathrm{D}\pi(\mathbf{P})[\mathbf{M}^{1/2}\mathbf{S}\mathbf{M}^{1/2}\mathbf{P}]]\mathrm{D}\pi(\mathbf{P})[\mathbf{M}^{1/2}\mathbf{H}\mathbf{M}^{1/2}\mathbf{P}]) \\
=&\frac{1}{2}\mathrm{tr}(\mathbf{S}\mathbf{M}\mathbf{H}\pi(\mathbf{P}) + \mathbf{S}\pi(\mathbf{P})\mathbf{H}\mathbf{M}) = \mathrm{tr}(\pi(\mathbf{P})\mathbf{S}\mathbf{M}\mathbf{H}),
\end{aligned}
$$

where the last equality is because $\mathbf{S}, \mathbf{M}, \mathbf{H}, \pi(\mathbf{P})$ are all symmetric. The inner product at $\mathbf{P}$ is given by

$$\langle \mathbf{M}^{1/2}\mathbf{S}\mathbf{M}^{1/2}\mathbf{P}, \mathbf{M}^{1/2}\mathbf{H}\mathbf{M}^{1/2}\mathbf{P}\rangle_2 = \operatorname{tr}(\mathbf{P}^\top\mathbf{M}^{1/2}\mathbf{S}\mathbf{M}\mathbf{H}\mathbf{M}^{1/2}\mathbf{P}) = \operatorname{tr}(\pi(\mathbf{P})\mathbf{S}\mathbf{M}\mathbf{H}).$$

This shows for any $\tilde{\mathbf{S}}, \tilde{\mathbf{H}} \in \mathcal{H}_{\mathbf{P}}$, $\langle \tilde{\mathbf{S}}, \tilde{\mathbf{H}}\rangle_2 = \langle \mathrm{D}\pi(\mathbf{P})[\tilde{\mathbf{S}}], \mathrm{D}\pi(\mathbf{P})[\tilde{\mathbf{H}}]\rangle_{\mathrm{gbw}}$, thereby completing the proof.

### B.5    Proof of Proposition 4

*Proof (Proof of Proposition 4).* Given for any $\mathbf{X} \in \mathbb{S}^n_{++}$, $\tau : \mathbb{S}^n_{++} \to \mathbb{S}^n_{++}$ is a diffeomorphism, it is thus suffices to show $g_{\mathrm{gbw},\mathbf{X}}(\mathbf{U}, \mathbf{V}) = g_{\mathrm{bw},\tau(\mathbf{X})}(\tau(\mathbf{U}), \tau(\mathbf{V}))$. That is,

$$\begin{aligned}
g_{\mathrm{gbw},\mathbf{X}}(\mathbf{U}, \mathbf{V}) &= \frac{1}{2}\operatorname{tr}(\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{V}) = \frac{1}{2}\operatorname{tr}(\mathbf{M}^{-1/2}\mathcal{L}_{\tau(\mathbf{X})}[\tau(\mathbf{U})]\mathbf{M}^{-1/2}\mathbf{V}) \\
&= \frac{1}{2}\operatorname{tr}(\mathcal{L}_{\tau(\mathbf{X})}[\tau(\mathbf{U})]\tau(\mathbf{V})) = g_{\mathrm{bw},\tau(\mathbf{X})}(\tau(\mathbf{U}), \tau(\mathbf{V})),
\end{aligned}$$

where we use the definition of the Lyapunov operator.

### B.6    Proof of Proposition 7

First we provide a Theorem that shows the pushforward distance from a Riemannian submersion is the Riemannian distance.

**Theorem 1 (Riemannian distance induced from Riemannian submersion [54]).** *Consider $\pi : (\mathcal{M}, g) \to (\mathcal{N}, h)$ as a Riemannian submersion. Let $d_{\mathcal{M}}$ be the Riemannian distance on $(\mathcal{M}, g)$ and the pushforward distance $d_{\mathcal{N}}(p, q) = \inf_{u \in \pi^{-1}(p), v \in \pi^{-1}(q)} d_{\mathcal{M}}(u, v)$ is equal to the Riemannian distance.*

We now proceed to derive the distance expression.

*Proof (Proof of Proposition 7).* From the definition of $\pi$ and Theorem 1, we have for any $\mathbf{X}, \mathbf{Y} \in \mathbb{S}^n_{++}$,

$$\begin{aligned}
d^2_{\mathrm{gbw}}(\mathbf{X}, \mathbf{Y}) &= \inf_{\mathbf{\Omega}, \mathbf{R} \in O(n)} d^2_{\mathrm{gl}}(\mathbf{M}^{-1/2}\mathbf{X}^{1/2}\mathbf{\Omega}, \mathbf{M}^{-1/2}\mathbf{Y}^{1/2}\mathbf{R}) \\
&= \inf_{\mathbf{\Omega}, \mathbf{R} \in O(n)} \|\mathbf{M}^{-1/2}\mathbf{X}^{1/2}\mathbf{\Omega} - \mathbf{M}^{-1/2}\mathbf{Y}^{1/2}\mathbf{R}\|^2_2 \\
&= \operatorname{tr}(\mathbf{M}^{-1}\mathbf{X}) + \operatorname{tr}(\mathbf{M}^{-1}\mathbf{Y}) - 2 \sup_{\mathbf{\Omega}, \mathbf{R} \in O(n)} \operatorname{tr}(\mathbf{M}^{-1}\mathbf{X}^{1/2}\mathbf{\Omega}\mathbf{R}^\top\mathbf{Y}^{1/2}) \\
&= \operatorname{tr}(\mathbf{M}^{-1}\mathbf{X}) + \operatorname{tr}(\mathbf{M}^{-1}\mathbf{Y}) - 2 \sup_{\mathbf{O} \in O(n)} \operatorname{tr}(\mathbf{M}^{-1}\mathbf{X}^{1/2}\mathbf{O}\mathbf{Y}^{1/2}).
\end{aligned}$$

The supremum is attained when $\mathbf{O} = \mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}^{1/2}(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{-1/2}$ as in Proposition 1. This completes the proof.

Here we also verify that the second-order approximation of the GBW distance recovers the proposed Riemannian metric in (3).

**Proposition 13.** *The GBW distance is approximated as* $d^2_{\mathrm{gbw}}(\mathbf{X}, \mathbf{X} + \theta\mathbf{H}) = \frac{\theta^2}{2}\mathrm{tr}(\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{H}]\mathbf{H}) + o(\theta^2)$.

*Proof (Proof of Proposition 13).* For $\mathbf{X} \in \mathbb{S}^n_{++}$ and $\mathbf{H} \in \mathbb{S}^n$ such that $\mathbf{X}\pm\mathbf{H} \in \mathbb{S}^n_{++}$. Thus, for $\theta \in [-1, 1]$, $\mathbf{X} + \theta\mathbf{H} \in \mathbb{S}^n_{++}$ and

$$
\begin{aligned}
d^2_{\mathrm{gbw}}(\mathbf{X}, \mathbf{X} + \theta\mathbf{H}) = {}& 2\mathrm{tr}(\mathbf{M}^{-1}\mathbf{X}) + \theta\mathrm{tr}(\mathbf{M}^{-1}\mathbf{H}) \\
& - 2\mathrm{tr}(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2} + \theta\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{1/2}
\end{aligned}
$$

The first-order derivative is

$$
\begin{aligned}
& \frac{d}{d\theta}d^2_{\mathrm{gbw}}(\mathbf{X}, \mathbf{X} + \theta\mathbf{H}) \\
& = \mathrm{tr}\Big(\mathbf{M}^{-1}\mathbf{H} \\
& \quad - 2\mathcal{L}_{(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2}+\theta\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{1/2}}[\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\mathbf{M}^{-1}\mathbf{X}^{1/2}]\Big) \\
& = \mathrm{tr}\Big(\mathbf{M}^{-1}\mathbf{H} \\
& \quad - (\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2} + \theta\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{-1/2}\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\mathbf{M}^{-1}\mathbf{X}^{1/2}\Big),
\end{aligned}
$$

where we use the properties of standard Lyapunov operator, $\mathbf{D}_{\mathbf{V}}(\mathbf{X})^{1/2} = \mathcal{L}_{\mathbf{X}^{1/2}}[\mathbf{V}]$ and $\mathrm{tr}(\mathcal{L}_{\mathbf{X}}[\mathbf{U}]) = \frac{1}{2}\mathrm{tr}(\mathbf{X}^{-1}\mathbf{U})$. Notice that

$$
\begin{aligned}
& \frac{d}{d\theta}d^2_{\mathrm{gbw}}(\mathbf{X}, \mathbf{X} + \theta\mathbf{H})|_{\theta=0} \\
& = \mathrm{tr}(\mathbf{M}^{-1}\mathbf{H}) - \mathrm{tr}\Big(\mathbf{M}^{-1}\mathbf{X}^{1/2}(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{-1/2}\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\Big) \\
& = \mathrm{tr}(\mathbf{M}^{-1}\mathbf{H}) - \mathrm{tr}\Big((\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X})^{-1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{H}\Big) = 0,
\end{aligned}
$$

where the second equality is from (11). The second order derivative is

$$
\begin{aligned}
& \frac{d^2}{d\theta^2}d^2_{\mathrm{gbw}}(\mathbf{X}, \mathbf{X} + \theta\mathbf{H}) \\
& = -\mathrm{tr}\Big(\frac{d}{d\theta}\big(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2} + \theta\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\mathbf{M}^{-1}\mathbf{X}^{1/2}\big)^{-1/2} \times \\
& \qquad\qquad \mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\mathbf{M}^{-1}\mathbf{X}^{1/2}\Big) \\
& = \mathrm{tr}(\frac{d}{d\theta}(-C^{-1/2})\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\mathbf{M}^{-1}\mathbf{X}^{1/2}),
\end{aligned}
$$

where we let $\mathbf{C} = \mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2} + \theta\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\mathbf{M}^{-1}\mathbf{X}^{1/2}$. Then,

$$
\frac{d}{d\theta}(-\mathbf{C}^{-1/2}) = \mathbf{C}^{-1/2}\mathcal{L}_{C^{1/2}}[\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\mathbf{M}^{-1}\mathbf{X}^{1/2}]\mathbf{C}^{-1/2}.
$$

Thus,

$$\frac{d^2}{d\theta^2}d^2_{\mathrm{gbw}}(\mathbf{X}, \mathbf{X} + \theta\mathbf{H})|_{\theta=0}$$
$$=\mathrm{tr}((\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{-1/2}\mathcal{L}_{(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{1/2}}[\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\mathbf{M}^{-1}\mathbf{X}^{1/2}]\times$$
$$(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{-1/2}\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\mathbf{M}^{-1}\mathbf{X}^{1/2})$$
$$=\mathrm{tr}(\mathbf{X}^{-1/2}\mathcal{L}_{(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{1/2}}[\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\mathbf{M}^{-1}\mathbf{X}^{1/2}]\mathbf{X}^{-1/2}\mathbf{H}).$$

Notice, similarly from (11),

$$(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{1/2}\mathbf{X}^{-1/2}\mathbf{M} = \mathbf{X}^{-1/2}\mathbf{M}(\mathbf{X}^{-1}\mathbf{M}\mathbf{X}^{-1}\mathbf{M})^{-1/2} = \mathbf{X}^{1/2}, \text{ and}$$
$$\mathbf{M}\mathbf{X}^{-1/2}(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{1/2} = \mathbf{X}^{1/2}$$

Let $\mathbf{L} := \mathcal{L}_{(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{1/2}}[\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{H}\mathbf{M}^{-1}\mathbf{X}^{1/2}]$. Then,

$$\begin{aligned}\mathbf{H} &= \mathbf{M}\mathbf{X}^{-1/2}\mathbf{L}(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{1/2}\mathbf{X}^{-1/2}\mathbf{M}\\ &\quad+ \mathbf{M}\mathbf{X}^{-1/2}(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{1/2}\mathbf{L}\mathbf{X}^{-1/2}\mathbf{M}\\ &= \mathbf{M}\mathbf{X}^{-1/2}\mathbf{L}\mathbf{X}^{1/2} + \mathbf{X}^{1/2}\mathbf{L}\mathbf{X}^{-1/2}\mathbf{M}\\ &= \mathbf{M}\mathbf{X}^{-1/2}\mathbf{L}\mathbf{X}^{-1/2}\mathbf{X} + \mathbf{X}\mathbf{X}^{-1/2}\mathbf{L}\mathbf{X}^{-1/2}\mathbf{M}.\end{aligned}$$

Thus, $\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{H}] = \mathbf{X}^{-1/2}\mathbf{L}\mathbf{X}^{-1/2}$ and $\frac{d^2}{d\theta^2}d^2_{\mathrm{gbw}}(\mathbf{X}, \mathbf{X} + \theta\mathbf{H})|_{\theta=0} = \mathrm{tr}(\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{H}]\mathbf{H})$. This completes the proof.

### B.7  An important lemma regarding the polar factor

The next lemma studies the various expressions of the polar factor $\mathbf{O}$, which is used throughout the proofs in the rest of the paper.

**Lemma 2.** *Consider* $\mathbf{O}$ *as defined in the proof of Proposition* (7)*, then*

$$\mathbf{O} = \mathbf{Y}^{1/2}(\mathbf{Y}^{-1}\mathbf{M}\mathbf{X}^{-1}\mathbf{M})^{1/2}\mathbf{M}^{-1}\mathbf{X}^{1/2} = \mathbf{Y}^{-1/2}(\mathbf{Y}\#(\mathbf{M}\mathbf{X}^{-1}\mathbf{M}))\mathbf{M}^{-1}\mathbf{X}^{1/2}.$$

*Proof.* From the definition of $\mathbf{O}$,

$$\begin{aligned}\mathbf{O} &= \mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}^{1/2}(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{-1/2}\\ &= \mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}^{1/2}(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{-1/2}\mathbf{X}^{-1/2}\mathbf{M}\mathbf{M}^{-1}\mathbf{X}^{1/2}\\ &= \mathbf{Y}^{1/2}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{Y})^{-1/2}\mathbf{M}^{-1}\mathbf{X}^{1/2} \quad\quad\quad\quad\quad (11)\\ &= \mathbf{Y}^{1/2}(\mathbf{Y}^{-1}\mathbf{M}\mathbf{X}^{-1}\mathbf{M})^{1/2}\mathbf{M}^{-1}\mathbf{X}^{1/2},\\ &= \mathbf{Y}^{-1/2}\mathbf{Y}\#(\mathbf{M}\mathbf{X}^{-1}\mathbf{M})\mathbf{M}^{-1}\mathbf{X}^{1/2}, \quad\quad\quad\quad\quad (12)\end{aligned}$$

where (11) is proved as follows. Denote $\mathbf{C} = (\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{-1/2}$ and we have

$$\begin{aligned}\mathbf{I} &= \mathbf{C}\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{X}^{1/2}\mathbf{C}\\ &= (\mathbf{M}^{-1}\mathbf{X}^{1/2}\mathbf{C}\mathbf{X}^{-1/2}\mathbf{M})\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{Y}(\mathbf{M}^{-1}\mathbf{X}^{1/2}\mathbf{C}\mathbf{X}^{-1/2}\mathbf{M}).\end{aligned}$$

Thus, $\mathbf{M}^{-1}\mathbf{X}^{1/2}\mathbf{C}\mathbf{X}^{-1/2}\mathbf{M} = (\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{Y})^{-1/2}$.

### B.8    Poof of Proposition 8

To derive the geodesic expression, we need the following well-known theorem.

**Theorem 2 (Geodesic induced from Riemannian submersion [8,36]).**
*Consider $\pi : (\mathcal{M}, g) \to (\mathcal{N}, h)$ as a Riemannian submersion. Let $c$ be a geodesic on $(\mathcal{M}, g)$ with $c'(0)$ is horizontal. Then, we have*

*(1) $c'(t)$ is horizontal for all $t$.*
*(2) $\gamma := \pi \circ c$ is a geodesic on $(\mathcal{N}, h)$ of the same length as $c$.*

*Proof (Proof of Proposition 8).* First, we see $\gamma(0) = \mathbf{X}, \gamma(1) = \mathbf{Y}$ and for $\mathbf{M}, \mathbf{X}, \mathbf{Y} \in \mathbb{S}_{++}^n$,

$$
\begin{aligned}
c(t) &= ((1-t)\mathbf{I} + t\mathbf{M}^{-1/2}\mathbf{Y}^{1/2}\mathbf{O}\mathbf{X}^{-1/2}\mathbf{M}^{1/2})\mathbf{M}^{-1/2}\mathbf{X}^{1/2} \\
&= ((1-t)\mathbf{I} + t\mathbf{M}^{-1/2}\mathbf{Y}\#(\mathbf{M}\mathbf{X}^{-1}\mathbf{M})\mathbf{M}^{-1/2})\mathbf{M}^{-1/2}\mathbf{X}^{1/2},
\end{aligned}
$$

where the second equality follows from Lemma 2. It is clear that

$$
\mathbf{M}^{-1/2}\mathbf{Y}\#(\mathbf{M}\mathbf{X}^{-1}\mathbf{M})\mathbf{M}^{-1/2} \in \mathbb{S}_{++}^n,
$$

and hence, $c(t)$ lies entirely in $\mathrm{GL}(n)$ for $t \in [0, 1]$ as it is closed under matrix multiplication. Also, $c(t)$ is a line segment, and thus, it is a valid geodesic on $\mathcal{M}_{\mathrm{gl}}$. Now, we need to show $c'(0)$ is horizontal. Indeed, we have

$$
\begin{aligned}
c'(0) &= \mathbf{M}^{-1/2}\mathbf{Y}^{1/2}\mathbf{O} - \mathbf{M}^{-1/2}\mathbf{X}^{1/2} \\
&= \mathbf{M}^{1/2}(\mathbf{M}^{-1}\mathbf{Y}^{1/2}\mathbf{O} - \mathbf{M}^{-1}\mathbf{X}^{1/2}) \\
&= \mathbf{M}^{1/2}(\mathbf{M}^{-1}\mathbf{Y}^{1/2}\mathbf{O}\mathbf{X}^{-1/2}\mathbf{M} - \mathbf{I})\mathbf{M}^{-1}\mathbf{X}^{1/2} \\
&= \mathbf{M}^{1/2}(\mathbf{M}^{-1}\mathbf{Y}(\mathbf{Y}^{-1}\mathbf{M}\mathbf{X}^{-1}\mathbf{M})^{1/2}\mathbf{M}^{-1} - \mathbf{M}^{-1})\mathbf{M}^{1/2}\mathbf{M}^{-1/2}\mathbf{X}^{1/2} \\
&= \mathbf{M}^{1/2}\mathbf{H}\mathbf{M}^{1/2}\mathbf{M}^{-1/2}\mathbf{X}^{1/2},
\end{aligned}
$$

where $\mathbf{H} := \mathbf{M}^{-1}\mathbf{Y}\#(\mathbf{M}\mathbf{X}^{-1}\mathbf{M})\mathbf{M}^{-1} - \mathbf{M}^{-1} \in \mathbb{S}^n$. Thus, from the definition of the horizontal space in (10), we have $c'(0) \in \mathcal{H}_{\mathbf{M}^{-1/2}\mathbf{X}^{1/2}}$. This completes the proof. In addition, from Theorem 2, we verify that the square of the Riemannian distance $d_{\mathrm{gbw}}^2$ is the same as the straight-line distance on $\mathcal{M}_{\mathrm{gl}}$, which is $\|\mathbf{M}^{-1/2}\mathbf{X}^{1/2} - \mathbf{M}^{-1/2}\mathbf{Y}^{1/2}\mathbf{O}\|_2^2 = \mathrm{tr}(\mathbf{M}^{-1}\mathbf{X}) + \mathrm{tr}(\mathbf{M}^{-1}\mathbf{Y}) - 2\mathrm{tr}(\mathbf{X}^{1/2}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{X}^{1/2})^{1/2}$.

### B.9    Proof of Proposition 9

*Proof (Proof of Proposition 9).* We first simplify $(1-t)\mathbf{X}^{1/2} + t\mathbf{Y}^{1/2}\mathbf{O} = ((1-t)\mathbf{I} + t\mathbf{Y}^{1/2}\mathbf{U}\mathbf{X}^{-1/2})\mathbf{X}^{1/2} = ((1-t)\mathbf{M} + t\mathbf{Y}\#(\mathbf{M}\mathbf{X}^{-1}\mathbf{M}))\mathbf{M}^{-1}\mathbf{X}^{1/2}$. With $\mathbf{K} := \mathbf{Y}\#(\mathbf{M}\mathbf{X}^{-1}\mathbf{M})$, we rewrite the geodesic as

$$
\begin{aligned}
\gamma(t) &= ((1-t)\mathbf{X}^{1/2} + t\mathbf{Y}^{1/2}\mathbf{O})((1-t)\mathbf{X}^{1/2} + t\mathbf{Y}^{1/2}\mathbf{O})^\top \\
&= ((1-t)\mathbf{M} + t\mathbf{K})\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}((1-t)\mathbf{M} + t\mathbf{K}) \\
&= \mathbf{X} + t\mathbf{X}(\mathbf{M}^{-1}\mathbf{K} - \mathbf{I}) + t(\mathbf{K}\mathbf{M}^{-1} - \mathbf{I})\mathbf{X} \\
&\quad + t^2\mathbf{M}(\mathbf{M}^{-1}\mathbf{K}\mathbf{M}^{-1} - \mathbf{M}^{-1})\mathbf{X}(\mathbf{M}^{-1}\mathbf{K}\mathbf{M}^{-1} - \mathbf{M}^{-1})\mathbf{M}.
\end{aligned}
$$

The first-order derivative is

$$\gamma'(0) = (\mathbf{K} - \mathbf{M})\mathbf{M}^{-1}\mathbf{X} + \mathbf{X}\mathbf{M}^{-1}(\mathbf{K} - \mathbf{M}) = (\mathbf{K}\mathbf{M}^{-1} - \mathbf{I})\mathbf{X} + \mathbf{X}(\mathbf{M}^{-1}\mathbf{K} - \mathbf{I})$$
$$= \mathbf{M}(\mathbf{M}^{-1}\mathbf{K}\mathbf{M}^{-1} - \mathbf{M}^{-1})\mathbf{X} + \mathbf{X}(\mathbf{M}^{-1}\mathbf{K}\mathbf{M}^{-1} - \mathbf{M}^{-1})\mathbf{M}.$$

Hence, $\gamma(t) = \mathbf{X} + t\gamma'(0) + t^2\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\gamma'(0)]\mathbf{X}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\gamma'(0)]\mathbf{M}$. The exponential map, therefore, is

$$\mathrm{Exp}_{\mathbf{X}}(t\mathbf{U}) = \mathbf{X} + t\mathbf{U} + t^2\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{X}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{M}$$
$$= (\mathbf{I} + t\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}])\mathbf{X}(\mathbf{I} + t\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{M})$$
$$= (\mathbf{M} + t\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{M})\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}(\mathbf{M} + t\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{M}).$$

Note that $\mathrm{Exp}_{\mathbf{X}}(t\mathbf{U}) \in \mathbb{S}_{++}^n$ if $\mathbf{M} + t\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{M} \in \mathbb{S}_{++}^n$.

To derive the logarithm map, let $\mathbf{Y} = \mathrm{Exp}_{\mathbf{X}}(\mathbf{U})$. We first have

$$\mathbf{M} + \mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{M}$$
$$= (\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{-1/2}\left((\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{1/2}\mathbf{Y}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{1/2}\right)^{1/2}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{-1/2}.$$

and

$$\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}] = -\mathbf{M}^{-1} + \mathbf{M}^{-1}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{-1/2}\times$$
$$\left((\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{1/2}\mathbf{Y}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{1/2}\right)^{1/2}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{-1/2}\mathbf{M}^{-1}$$

Hence, let $\mathbf{S}_{\mathbf{M}} := \left((\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{1/2}\mathbf{Y}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{1/2}\right)^{1/2}$. Then,

$$\mathbf{U} = \mathbf{X}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{M} + \mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{X}$$
$$= 2\{\mathbf{X}\mathbf{M}^{-1}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{-1/2}\mathbf{S}_{\mathbf{M}}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{-1/2}\}_{\mathrm{S}} - 2\mathbf{X}$$
$$= 2\{\mathbf{M}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{1/2}\mathbf{S}_{\mathbf{M}}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{-1/2}\}_{\mathrm{S}} - 2\mathbf{X}$$
$$= \mathbf{M}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{Y})^{1/2} + (\mathbf{Y}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{1/2}\mathbf{M} - 2\mathbf{X},$$

where we denote $\{\mathbf{A}\}_{\mathrm{S}} := (\mathbf{A} + \mathbf{A}^\top)/2$, for $\mathbf{A} \in \mathbb{R}^{n\times n}$. This completes the proof.

### B.10   Proof of Proposition 10

The Levi-Civita connection (or Levi-Civita derivative) of a vector field on a manifold $\mathcal{M}$ is the unique covariant derivative that satisfies (1) torsion-free property, i.e., $\nabla_\xi \eta - \nabla_\eta \xi = \mathrm{D}_\xi \eta - \mathrm{D}_\eta \xi = [\xi, \eta]$ and (2) metric compatibility, i.e., $\nabla_\xi \langle \eta, \xi \rangle_{\mathcal{M}} = \langle \nabla_\xi \eta, \zeta \rangle_{\mathcal{M}} + \langle \eta, \nabla_\xi \zeta \rangle_{\mathcal{M}}$, for any vector fields $\xi, \eta, \zeta$.

*Proof (Proof of Proposition 10).* The Levi-Civita connection is derived by applying [33, MD.3]. For any vector fields $\xi, \eta, \zeta$ on $\mathcal{M}_{\mathrm{gbw}}$, it satisfies for any $\mathbf{X} \in \mathcal{M}_{\mathrm{gbw}}$,

$$\langle \nabla_\xi \eta, \mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta] \rangle_2$$

$$= \langle \mathrm{D}_\xi \eta, \mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta] \rangle_2 + \frac{1}{2}\langle \eta, \mathrm{D}_\xi \mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta] \rangle_2 + \frac{1}{2}\langle \xi, \mathrm{D}_\eta \mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta] \rangle_2 - \frac{1}{2}\langle \xi, \mathrm{D}_\zeta \mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta] \rangle$$

$$= \langle \mathrm{D}_\xi \eta, \mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta] \rangle_2 + \frac{1}{2}\langle \xi, \mathcal{L}_{\mathbf{X},\mathbf{M}}\left[\zeta \mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\mathbf{M} + \mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\zeta\right]\rangle_2$$

$$\quad - \frac{1}{2}\langle \eta, \mathcal{L}_{\mathbf{X},\mathbf{M}}\left[\xi \mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta]\mathbf{M} + \mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta]\xi\right]\rangle_2$$

$$\quad - \frac{1}{2}\langle \xi, \mathcal{L}_{\mathbf{X},\mathbf{M}}\left[\eta \mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta]\mathbf{M} + \mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta]\eta\right]\rangle_2. \tag{13}$$

The second term of (13) is rewritten as

$$\frac{1}{2}\langle \xi, \mathcal{L}_{\mathbf{X},\mathbf{M}}\left[\zeta \mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\mathbf{M} + \mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\zeta\right]\rangle_2$$

$$= \frac{1}{2}\langle \mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi] + \mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta], \zeta \rangle_2$$

$$= \langle \{\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\}_{\mathrm{S}}, \zeta \rangle_2$$

$$= \langle \mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{X}\{\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\}_{\mathrm{S}}\mathbf{M} + \mathbf{M}\{\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\}_{\mathrm{S}}\mathbf{X}], \zeta \rangle_2$$

$$= \langle \mathbf{X}\{\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\}_{\mathrm{S}}\mathbf{M} + \mathbf{M}\{\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\}_{\mathrm{S}}\mathbf{X}, \mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta] \rangle_2$$

$$= \langle \{\mathbf{X}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\mathbf{M} + \mathbf{X}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\mathbf{M}\}_{\mathrm{S}}, \mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta] \rangle_2. \tag{14}$$

Similarly,

$$\frac{1}{2}\langle \eta, \mathcal{L}_{\mathbf{X},\mathbf{M}}\left[\xi \mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta]\mathbf{M} + \mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta]\xi\right]\rangle_2 = \langle \{\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\eta]\xi\}_{\mathrm{S}}, \mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta] \rangle_2 \tag{15}$$

$$\frac{1}{2}\langle \xi, \mathcal{L}_{\mathbf{X},\mathbf{M}}\left[\eta \mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta]\mathbf{M} + \mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta]\eta\right]\rangle_2 = \langle \{\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\xi]\eta\}_{\mathrm{S}}, \mathcal{L}_{\mathbf{X},\mathbf{M}}[\zeta] \rangle_2. \tag{16}$$

Applying the results in (14), (15), and (16) in (13), the proof is complete.

### B.11   Proof of Proposition 11

We first provide the formal definition of sectional curvature. The curvature tensor $R$ is defined for any $X, Y, Z \in \mathfrak{X}(\mathcal{M})$, $R(X,Y)Z := \nabla_X \nabla_Y Z - \nabla_Y \nabla_X Z - \nabla_{[X,Y]}Z$, where $[X,Y] = XY - YX$ is the Lie bracket and $\nabla$ is the Levi-Civita connection. At point $p$, $R_p$ defines a $(1,3)$-tensor on $T_p\mathcal{M}$ with $R_p(X(p), Y(p))Z(p) \in T_p\mathcal{M}$, where $X(p), Y(p), Z(p) \in T_p\mathcal{M}$ are the vector fields evaluated at $p \in \mathcal{M}$. The sectional curvature is the scalar curvature of a 2-dimensional subspace of $T_p\mathcal{M}$, given by

$$K(u,v) = \frac{g(R_p(u,v)v, u)}{g(u,u)g(v,v) - (g(u,v))^2} \tag{17}$$

for $u, v \in T_p\mathcal{M}$ as two linearly independent tangent vectors that span the subspace.

Before deriving the sectional curvature, we require the following theorem from Riemannian submersion.

We now derive the sectional curvature of $\mathcal{M}_{\mathrm{gbw}}$ based on the following theorem from Riemannian submersion [6,44].

**Theorem 3.** *Let $\pi : (\tilde{\mathcal{M}}, \tilde{g}) \to (\mathcal{M}, g)$ be a Riemannian submersion and consider $X, Y$ as smooth vector fields on $\mathcal{M}$. The horizontal lift $\tilde{X}, \tilde{Y}$ are unique vector fields on $\tilde{\mathcal{M}}$ such that $\tilde{X}(p), \tilde{Y}(p) \in \mathcal{H}_p$ and $\mathrm{D}\pi(p)[\tilde{X}(p)] = X(\pi(p)), \mathrm{D}\pi(p)[\tilde{Y}(p)] = Y(\pi(p))$ for all $p \in \tilde{\mathcal{M}}$. Then, the sectional curvature is*

$$K(X, Y) = \tilde{K}(\tilde{X}, \tilde{Y}) + \frac{3}{4} \frac{\|[\tilde{X}, \tilde{Y}]^{\mathcal{V}}\|^2}{Q(\tilde{X}, \tilde{Y})},$$

*where $Q(\tilde{X}, \tilde{Y}) = \tilde{g}(\tilde{X}, \tilde{X})\tilde{g}(\tilde{Y}, \tilde{Y}) - (\tilde{g}(\tilde{X}, \tilde{Y}))^2$. $Z^{\mathcal{V}}$ is the vertical component of a vector field and $\tilde{K}$ is the sectional curvature of $(\tilde{\mathcal{M}}, \tilde{g})$.*

Directly from Theorem 3 above, we see that the sectional curvature of $\mathcal{M}_{\mathrm{gbw}}$ is non-negative, given that $\mathcal{M}_{\mathrm{gl}}$ endowed with the flat Euclidean metric has zero curvature. Before we derive the sectional curvature, we need the following lemma to show projection to the horizontal/vertical space on $T_{\mathbf{P}}\mathcal{M}_{\mathrm{gl}}$.

**Lemma 3.** *Any $\mathbf{U} \in T_{\mathbf{P}}\mathcal{M}_{\mathrm{gl}}$ can be projected onto the vertical and horizontal spaces defined in Proposition 3, i.e., $\mathbf{U} = \mathbf{U}^{\mathcal{V}} + \mathbf{U}^{\mathcal{H}}$, where*

$$\mathbf{U}^{\mathcal{H}} = \mathbf{M}^{1/2}\mathcal{L}_{\mathbf{M},\mathbf{M}^{1/2}\mathbf{PP}^{\top}\mathbf{M}^{1/2}}[\mathbf{M}^{1/2}(\mathbf{UP}^{\top} + \mathbf{PU}^{\top})\mathbf{M}^{1/2}]\mathbf{M}^{1/2}\mathbf{P},$$
$$\mathbf{U}^{\mathcal{V}} = \mathbf{M}^{-1/2}\mathcal{L}_{\mathbf{M}^{-1},(\mathbf{M}^{1/2}\mathbf{PP}^{\top}\mathbf{M}^{1/2})^{-1}}[\mathbf{M}^{-1/2}(\mathbf{UP}^{-1} - \mathbf{P}^{-\top}\mathbf{U}^{\top})\mathbf{M}^{-1/2}]\mathbf{M}^{-1/2}\mathbf{P}^{-\top}.$$

*Proof.* Based on Proposition 3, for $\mathbf{U} \in T_{\mathbf{P}}\mathcal{M}_{\mathrm{gl}}$, it can be decomposed as $\mathbf{U} = \mathbf{U}^{\mathcal{V}} + \mathbf{U}^{\mathcal{H}} = \mathbf{M}^{-1/2}\mathbf{KM}^{-1/2}\mathbf{P}^{-\top} + \mathbf{M}^{1/2}\mathbf{SM}^{1/2}\mathbf{P}$, for $\mathbf{K}$ skew-symmetric and $\mathbf{S}$ symmetric. From the decomposition, $\mathbf{U}^{\top} = -\mathbf{P}^{-1}\mathbf{M}^{-1/2}\mathbf{KM}^{-1/2} + \mathbf{P}^{\top}\mathbf{M}^{1/2}\mathbf{SM}^{1/2}$. Thus, we have

$$\mathbf{M}^{1/2}(\mathbf{UP}^{\top} + \mathbf{PU}^{\top})\mathbf{M}^{1/2} = \mathbf{MSM}^{1/2}\mathbf{PP}^{\top}\mathbf{M}^{1/2} + \mathbf{M}^{1/2}\mathbf{PP}^{\top}\mathbf{M}^{1/2}\mathbf{SM}.$$

Hence, $\mathbf{S} = \mathcal{L}_{\mathbf{M},\mathbf{M}^{1/2}\mathbf{PP}^{\top}\mathbf{M}^{1/2}}[\mathbf{M}^{1/2}(\mathbf{UP}^{\top} + \mathbf{PU}^{\top})\mathbf{M}^{1/2}]$. Similarly, we also have

$$\mathbf{M}^{-1/2}(\mathbf{UP}^{-1} - \mathbf{P}^{-\top}\mathbf{U}^{\top})\mathbf{M}^{-1/2}$$
$$= \mathbf{M}^{-1}\mathbf{KM}^{-1/2}\mathbf{P}^{-\top}\mathbf{P}^{-1}\mathbf{M}^{-1/2} + \mathbf{M}^{-1/2}\mathbf{P}^{-\top}\mathbf{P}^{-1}\mathbf{M}^{-1/2}\mathbf{KM}^{-1}$$

Thus, $\mathbf{K} = \mathcal{L}_{\mathbf{M}^{-1},(\mathbf{M}^{1/2}\mathbf{PP}^{\top}\mathbf{M}^{1/2})^{-1}}[\mathbf{M}^{-1/2}(\mathbf{UP}^{-1} - \mathbf{P}^{-\top}\mathbf{U}^{\top})\mathbf{M}^{-1/2}]$, which is clearly skew-symmetric given that $\mathbf{UP}^{-1} - \mathbf{P}^{-\top}\mathbf{U}^{\top}$ is skew-symmetric.

Finally, we proceed to prove the main proposition.

*Proof (Proof of Proposition 11).* We denote $\mathbf{S}_U := \mathcal{L}_{\mathbf{M},\pi(\mathbf{P})}[U(\pi(\mathbf{P}))]$ and similarly for $\mathbf{S}_V$. Hence we see $\tilde{U}, \tilde{V}$ are the horizontal lift according to the definition.

To start, it is clear $\tilde{U}(\mathbf{P}) \in \mathcal{H}_{\mathbf{P}}$ according to the definition of the horizontal space in Proposition 3. Also, we have

$$\mathrm{D}\pi(\mathbf{P})[\tilde{U}(\mathbf{P})]$$
$$= \mathbf{M}\mathcal{L}_{\mathbf{M},\pi(\mathbf{P})}[U(\pi(\mathbf{P}))]\mathbf{M}^{1/2}\mathbf{P}\mathbf{P}^{\top}\mathbf{M}^{1/2} + \mathbf{M}^{1/2}\mathbf{P}\mathbf{P}^{\top}\mathbf{M}^{1/2}\mathcal{L}_{\mathbf{M},\pi(\mathbf{P})}[U(\pi(\mathbf{P}))]\mathbf{M}$$
$$= U(\pi(\mathbf{P})), \qquad \forall \mathbf{P} \in \mathcal{M}_{\mathrm{gl}}.$$

This suggests $\tilde{U} \in \mathfrak{X}(\mathcal{M}_{\mathrm{gl}})$ is indeed a horizontal lift of $U \in \mathfrak{X}(\mathcal{M}_{\mathrm{gbw}})$. Next we compute the sectional curvature following Theorem 3.

First, we derive an expression for the Lie bracket. For any two horizontal tangent vectors $\tilde{U}(\mathbf{P}), \tilde{V}(\mathbf{P})$, they can be written as $\tilde{U}(\mathbf{P}) = \mathbf{M}^{1/2}\mathbf{S}_U\mathbf{M}^{1/2}\mathbf{P}$ and $\tilde{V}(\mathbf{P}) = \mathbf{M}^{1/2}\mathbf{S}_V\mathbf{M}^{1/2}\mathbf{P}$, for arbitrary symmetric matrices $\mathbf{S}_U, \mathbf{S}_V$. Therefore,

$$[\tilde{U}, \tilde{V}](\mathbf{P})$$
$$= \mathrm{D}\tilde{V}(\mathbf{P})[\tilde{U}(\mathbf{P})] - \mathrm{D}\tilde{U}(\mathbf{P})[\tilde{V}(\mathbf{P})]$$
$$= \mathbf{M}^{1/2}\mathrm{D}\mathbf{S}_V[\tilde{U}(\mathbf{P})]\mathbf{M}^{1/2}\mathbf{P} + \mathbf{M}^{1/2}\mathbf{S}_V\mathbf{M}^{1/2}\tilde{U}(\mathbf{P}) - \mathbf{M}^{1/2}\mathrm{D}\mathbf{S}_U[\tilde{V}(\mathbf{P})]\mathbf{M}^{1/2}\mathbf{P}$$
$$- \mathbf{M}^{1/2}\mathbf{S}_U\mathbf{M}^{1/2}\tilde{V}(\mathbf{P}).$$

From Lemma 3, to project the result onto the vertical space, we need to first evaluate

$$\mathbf{M}^{-1/2}\big(([\tilde{U}, \tilde{V}](\mathbf{P}))\mathbf{P}^{-1} - \mathbf{P}^{-\top}([\tilde{U}, \tilde{V}](\mathbf{P}))^{\top}\big)\mathbf{M}^{-1/2}$$
$$= \mathrm{D}\mathbf{S}_V[\tilde{U}(\mathbf{P})] + \mathbf{S}_V\mathbf{M}\mathbf{S}_U - \mathrm{D}\mathbf{S}_U[\tilde{V}(\mathbf{P})] - \mathbf{S}_U\mathbf{M}\mathbf{S}_V - \mathrm{D}\mathbf{S}_V[\tilde{U}(\mathbf{P})] - \mathbf{S}_U\mathbf{M}\mathbf{S}_V$$
$$+ \mathrm{D}\mathbf{S}_V[\tilde{U}(\mathbf{P})] + \mathbf{S}_V\mathbf{M}\mathbf{S}_U$$
$$= 2(\mathbf{S}_V\mathbf{M}\mathbf{S}_U - \mathbf{S}_U\mathbf{M}\mathbf{S}_V),$$

and the vertical projection is

$$([\tilde{U}, \tilde{V}](\mathbf{P}))^{\mathcal{V}} = \mathbf{M}^{-1/2}\mathcal{L}_{\mathbf{M}^{-1},\pi(\mathbf{P})^{-1}}[2\mathbf{S}_V\mathbf{M}\mathbf{S}_U - 2\mathbf{S}_U\mathbf{M}\mathbf{S}_V]\mathbf{M}^{-1/2}\mathbf{P}^{-\top}.$$

To study the trace norm of the vertical projection, we denote

$$\mathbf{L} := \mathcal{L}_{\mathbf{M}^{-1},\pi(\mathbf{P})^{-1}}[2\mathbf{S}_V\mathbf{M}\mathbf{S}_U - 2\mathbf{S}_U\mathbf{M}\mathbf{S}_V].$$

Then, from the definition of generalized Lyapunov operator,

$$\mathbf{P}^{\top}\mathbf{M}^{-1/2}\mathbf{L}\mathbf{M}^{-1/2}\mathbf{P}^{-\top} + \mathbf{P}^{-1}\mathbf{M}^{-1/2}\mathbf{L}\mathbf{M}^{-1/2}\mathbf{P}$$
$$= 2\mathbf{P}^{\top}\mathbf{M}^{1/2}(\mathbf{S}_V\mathbf{M}\mathbf{S}_U - \mathbf{S}_U\mathbf{M}\mathbf{S}_V)\mathbf{M}^{1/2}\mathbf{P}$$
$$= 2\tilde{V}(\mathbf{P})^{\top}\tilde{U}(\mathbf{P}) - 2\tilde{U}(\mathbf{P})^{\top}\tilde{V}(\mathbf{P}).$$

Now, consider the singular value decomposition of $\mathbf{P} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^{\top}$ with the singular values sorted decreasingly. Denote $\mathbf{C} := \mathbf{V}^{\top}(\tilde{V}(\mathbf{P})^{\top}\tilde{U}(\mathbf{P}) - \tilde{U}(\mathbf{P})^{\top}\tilde{V}(\mathbf{P}))\mathbf{V}$. This yields

$$2\mathbf{C} = \mathbf{\Sigma}\mathbf{U}^{\top}\mathbf{M}^{-1/2}\mathbf{L}\mathbf{M}^{-1/2}\mathbf{U}\mathbf{\Sigma}^{-1} + \mathbf{\Sigma}^{-1}\mathbf{U}^{\top}\mathbf{M}^{-1/2}\mathbf{L}\mathbf{M}^{-1/2}\mathbf{U}\mathbf{\Sigma}. \qquad (18)$$

Denote $\tilde{\mathbf{L}} := \mathbf{U}^\top \mathbf{M}^{-1/2} \mathbf{L} \mathbf{M}^{-1/2} \mathbf{U}$. Result (18) indicates $(\sigma_i \sigma_j^{-1} + \sigma_i^{-1} \sigma_j) \tilde{\mathbf{L}}_{ij} = 2\mathbf{C}_{ij}$. Hence, $\mathbf{M}^{-1/2} \mathbf{L} \mathbf{M}^{-1/2} = \mathbf{U} \tilde{\mathbf{L}} \mathbf{U}^\top$ and

$$
\begin{aligned}
\|([\tilde{U}, \tilde{V}](\mathbf{P}))^{\mathcal{V}}\|_2^2 = \|\mathbf{U} \tilde{\mathbf{L}} \mathbf{U}^\top \mathbf{P}^{-\top}\|_2^2 &= \|\mathbf{U} \tilde{\mathbf{L}} \boldsymbol{\Sigma}^{-1} \mathbf{V}^\top\|_2^2 \\
&= \|\tilde{\mathbf{L}} \boldsymbol{\Sigma}^{-1}\|_2^2 \\
&= \sum_{i,j} \frac{4 \mathbf{C}_{ij}^2}{\sigma_j^2 (\sigma_i \sigma_j^{-1} + \sigma_i^{-1} \sigma_j)^2}.
\end{aligned}
$$

Based on Theorem 3, the proof is complete by noticing $\mathcal{M}_{\mathrm{gl}}$ has zero curvature and choosing orthonormal tangent vectors $\tilde{U}(\mathbf{P}), \tilde{V}(\mathbf{P})$ without loss of generality.

### B.12  Proof of Proposition 12

We now compute the bounds for the sectional curvature following [42]. We need the following lemma, which bounds the skew operation of matrix product.

**Lemma 4 (Lemma 2 in [42]).** *For arbitrary matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ with $\|\mathbf{A}\|_2 = \|\mathbf{B}\|_2 = 1$, we have $\|\mathbf{A}^\top \mathbf{B} - \mathbf{B}^\top \mathbf{A}\|_2^2 \leq 2$.*

*Proof (Proof of Proposition 12).* It is clear when $\mathbf{C} = \mathbf{0}$, the sectional curvature is zero, which happens when for example, $\mathbf{S}_U = \mathbf{M}^{-1}, \mathbf{S}_V = \mathbf{S}$ for arbitrary symmetric $\mathbf{S}$. This holds even when $\tilde{U}(\mathbf{P}), \tilde{V}(\mathbf{P})$ are not orthonormal.

Also, we have

$$
\begin{aligned}
&K(U(\pi(\mathbf{P})), V(\pi(\mathbf{P}))) \\
&= \sum_{i,j} \frac{3\mathbf{C}_{ij}^2}{\sigma_j^2 (\sigma_i \sigma_j^{-1} + \sigma_i^{-1} \sigma_j)^2} = \sum_{i,j} \frac{3\sigma_i^2 \mathbf{C}_{ij}^2}{(\sigma_i^2 + \sigma_j^2)^2} = \frac{3 \sum_{i>j} (\sigma_i^2 + \sigma_j^2) \mathbf{C}_{ij}^2}{(\sigma_i^2 + \sigma_j^2)^2} \\
&= \sum_{i>j} \frac{3\mathbf{C}_{ij}^2}{\sigma_i^2 + \sigma_j^2} \leq \frac{3}{2(\sigma_n^2 + \sigma_{n-1}^2)} \|\mathbf{C}\|_2^2 \leq \frac{3}{\sigma_n^2 + \sigma_{n-1}^2},
\end{aligned}
$$

where we notice $\mathbf{C}$ is skew-symmetric and apply Lemma 4. To verify the choice of $\tilde{U}(\mathbf{P}), \tilde{V}(\mathbf{P})$ that achieves the maximum curvature, we first see

$$
\begin{aligned}
\operatorname{tr}(\tilde{U}(\mathbf{P})^\top \tilde{V}(\mathbf{P})) &= \frac{\operatorname{tr}(\boldsymbol{\Sigma}^{-2}(\mathbf{E}_{\{n-1,n-1\}} - \mathbf{E}_{\{n,n\}}) \mathbf{E}_{\{n,n-1\}})}{2(\sigma_n^{-2} + \sigma_{n-1}^{-2})} \\
&= \frac{\operatorname{tr}(\boldsymbol{\Sigma}^{-2}(\mathbf{e}_{n-1} \mathbf{e}_n^\top - \mathbf{e}_n \mathbf{e}_{n-1}^\top))}{\sigma_n^{-2} + \sigma_{n-1}^{-2}} = 0, \\
\operatorname{tr}(\tilde{U}(\mathbf{P})^\top \tilde{U}(\mathbf{P})) = \operatorname{tr}(\tilde{V}(\mathbf{P})^\top \tilde{V}(\mathbf{P})) &= \frac{\operatorname{tr}(\boldsymbol{\Sigma}^{-2}(\mathbf{e}_n \mathbf{e}_n^\top + \mathbf{e}_{n-1} \mathbf{e}_{n-1}^\top))}{\sigma_n^{-2} + \sigma_{n-1}^{-2}} = 1,
\end{aligned}
$$

which shows $\tilde{U}(\mathbf{P}), \tilde{V}(\mathbf{P})$ are orthonormal.

Also, we have

$$
\begin{aligned}
\mathbf{C} &= \mathbf{V}^\top(\tilde{V}(\mathbf{P})^\top\tilde{U}(\mathbf{P}) - \tilde{U}(\mathbf{P})^\top\tilde{V}(\mathbf{P}))\mathbf{V} \\
&= \frac{\mathbf{E}_{\{n,n-1\}}\boldsymbol{\Sigma}^{-2}(\mathbf{E}_{\{n-1,n-1\}} - \mathbf{E}_{\{n,n\}}) - (\mathbf{E}_{\{n-1,n-1\}} - \mathbf{E}_{\{n,n\}})\boldsymbol{\Sigma}^{-2}\mathbf{E}_{\{n,n-1\}}}{2(\sigma_n^{-2} + \sigma_{n-1}^{-2})} \\
&= \mathbf{e}_n\mathbf{e}_{n-1}^\top - \mathbf{e}_{n-1}\mathbf{e}_n^\top.
\end{aligned}
$$

This leads to the maximum sectional curvature as $\sum_{i>j}\frac{3\mathbf{C}_{ij}^2}{\sigma_i^2+\sigma_j^2} = \frac{3}{\sigma_n^2+\sigma_{n-1}^2}$.

### B.13    Proof of Proposition 15

*Proof (Proof of Proposition 15).* From the expression of GBW geodesic, we have

$$
\begin{aligned}
\gamma(t) &= (1-t)^2\mathbf{X} + t^2\mathbf{Y} + t(1-t)\left((\mathbf{Y}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1})^{1/2}\mathbf{M} + \mathbf{M}(\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{Y})^{1/2}\right) \\
&= (1-t)^2\mathbf{X} + t^2\mathbf{Y} + t(1-t)\Big(\mathbf{Y}^{1/2}(\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{Y}^{1/2})^{1/2}\mathbf{Y}^{-1/2}\mathbf{M} \\
&\quad + \mathbf{M}\mathbf{Y}^{-1/2}(\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{Y}^{1/2})^{1/2}\mathbf{Y}^{1/2}\Big) \\
&= \mathbf{M}\mathbf{Y}^{-1/2}\Big((1-t)^2(\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{Y}^{1/2}) + t^2(\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{Y}^{1/2}) \\
&\quad + t(1-t)\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{Y}^{1/2}(\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{Y}^{1/2})^{1/2} \\
&\quad + t(1-t)(\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{Y}^{1/2})^{1/2}\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{Y}^{1/2}\Big)\mathbf{Y}^{-1/2}\mathbf{M} \\
&= \mathbf{M}\mathbf{Y}^{-1/2}\Big((1-t)(\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{Y}^{1/2})^{1/2} + t(\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{Y}^{1/2})\Big)^2\mathbf{Y}^{-1/2}\mathbf{M} \\
&\preceq \mathbf{M}\mathbf{Y}^{-1/2}\Big((1-t)(\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{X}\mathbf{M}^{-1}\mathbf{Y}^{1/2}) + t(\mathbf{Y}^{1/2}\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}\mathbf{Y}^{1/2})\Big)\mathbf{Y}^{-1/2}\mathbf{M} \\
&= (1-t)\mathbf{X} + t\mathbf{Y},
\end{aligned}
$$

where the second equality follows from the property of geometric mean $(\mathbf{A}\mathbf{B})^{1/2} = \mathbf{A}(\mathbf{A}^{-1}\mathbf{B})^{1/2} = \mathbf{A}^{1/2}(\mathbf{A}^{1/2}\mathbf{B}\mathbf{A}^{1/2})^{1/2}\mathbf{A}^{-1/2}$.

### B.14    Proof of Theorem 4

*Proof (Proof of Theorem 4).* First we see

$$
F(\mathbf{A}) = \sum_{l=1}^{N} w_l\mathrm{tr}(\mathbf{M}^{-1}\mathbf{X}_l) + \sum_{l=1}^{N} w_l\mathrm{tr}\Big(\mathbf{M}^{-1}\mathbf{A} - 2(\mathbf{X}_l^{1/2}\mathbf{M}^{-1}\mathbf{A}\mathbf{M}^{-1}\mathbf{X}_l^{1/2})^{1/2}\Big).
$$

Thus to show strict convexity of $F(\mathbf{A})$, we only need to show

$$
S(\mathbf{A}) = \mathrm{tr}(\mathbf{X}_l^{1/2}\mathbf{M}^{-1}\mathbf{A}\mathbf{M}^{-1}\mathbf{X}_l^{1/2})^{1/2}
$$

is strictly concave. This is true because $\mathrm{tr}(\mathbf{X})^{1/2}$ is strictly concave. See proof in [8,5].

By first-order stationarity, we need to find the derivative of $F(\mathbf{A})$. First we write $S(\mathbf{A}) = \mathrm{tr}((h \circ \phi)(\mathbf{A}))$, where $h(\mathbf{A}) = \mathbf{A}^{1/2}$ and $\phi(\mathbf{A}) = \mathbf{X}_l^{1/2}\mathbf{M}^{-1}\mathbf{A}\mathbf{M}^{-1}\mathbf{X}_l^{1/2}$. Recall that $\mathrm{D}h(\mathbf{X})[\mathbf{U}] = \mathcal{L}_{\mathbf{X}^{1/2}}[\mathbf{U}]$ by the derivative of the inverse function law [41,8]. Thus by chain rule,

$$
\begin{aligned}
\mathrm{D}S(\mathbf{A})[\mathbf{U}] &= \mathrm{tr}\Big((\mathrm{D}h(\phi(\mathbf{A})) \circ \mathrm{D}\phi(\mathbf{A}))[\mathbf{U}]\Big) \\
&= \mathrm{tr}\Big(\mathcal{L}_{(\mathbf{X}_l^{1/2}\mathbf{M}^{-1}\mathbf{A}\mathbf{M}^{-1}\mathbf{X}_l^{1/2})^{1/2}}[\mathbf{X}_l^{1/2}\mathbf{M}^{-1}\mathbf{U}\mathbf{M}^{-1}\mathbf{X}_l^{1/2}]\Big) \\
&= \frac{1}{2}\mathrm{tr}\Big(\mathbf{M}^{-1}\mathbf{X}_l^{1/2}(\mathbf{X}_l^{1/2}\mathbf{M}^{-1}\mathbf{A}\mathbf{M}^{-1}\mathbf{X}_l^{1/2})^{-1/2}\mathbf{X}_l^{1/2}\mathbf{M}^{-1}\mathbf{U}\Big) \\
&= \frac{1}{2}\mathrm{tr}\Big((\mathbf{A}^{-1}\mathbf{M}\mathbf{X}_l^{-1}\mathbf{M})^{1/2}\mathbf{M}^{-1}\mathbf{X}_l\mathbf{M}^{-1}\mathbf{U}\Big) \\
&= \frac{1}{2}\mathrm{tr}\Big(\mathbf{A}^{-1}\#(\mathbf{M}^{-1}\mathbf{X}_l\mathbf{M}^{-1})\mathbf{U}\Big),
\end{aligned}
$$

where the second equality follows from $\mathrm{tr}(\mathcal{L}_{\mathbf{X}}[\mathbf{U}]) = \frac{1}{2}\mathrm{tr}(\mathbf{X}^{-1}\mathcal{L}_{\mathbf{X}}[\mathbf{U}]\mathbf{X} + \mathcal{L}_{\mathbf{X}}[\mathbf{U}]) = \frac{1}{2}\mathrm{tr}(\mathbf{X}^{-1}\mathbf{U})$ and the third equality is due to (12).

Hence, $\mathrm{D}F(\mathbf{A})[\mathbf{U}] = \sum_l w_l \mathrm{tr}\Big(\mathbf{M}^{-1}\mathbf{U} - \mathbf{A}^{-1}\#(\mathbf{M}^{-1}\mathbf{X}_l\mathbf{M}^{-1})\mathbf{U}\Big)$. From the first-order optimality of the convex function $F(\mathbf{A})$, i.e. $\mathrm{D}F(\mathbf{A})[\mathbf{U}] = \mathbf{0}$ for all $\mathbf{U}$, the unique minimizer $A(\mathbf{X}_{1:N}, \mathbf{w})$ satisfies $\mathbf{M}^{-1} = \sum_{l=1}^{N} w_l \mathbf{A}^{-1}\#(\mathbf{M}^{-1}\mathbf{X}_l\mathbf{M}^{-1})$, which is equivalent to $\mathbf{A}^{1/2}\mathbf{M}^{-1}\mathbf{A}^{1/2} = \sum_{l=1}^{N} w_l (\mathbf{A}^{1/2}\mathbf{M}^{-1}\mathbf{X}_l\mathbf{M}^{-1}\mathbf{A}^{1/2})^{1/2}$.

### B.15   Proof of Theorem 5

*Proof (Proof of Theorem 5).* First by convexity of the matrix square,

$$
K(\mathbf{A}) \leq \mathbf{M}\mathbf{A}^{-1/2}(\sum_{l=1}^{N} w_l \mathbf{A}^{1/2}\mathbf{M}^{-1}\mathbf{X}_l\mathbf{M}^{-1}\mathbf{A}^{1/2})\mathbf{A}^{-1/2}\mathbf{M} \leq \sum_{l=1}^{N}\mathbf{X}_l.
$$

Hence, $K(\mathbf{A})$ is bounded in $\mathbb{S}_{++}^n$. Also, we claim $F(\mathbf{A}_{t+1}) \leq F(\mathbf{A}_t)$, where $F(\mathbf{A})$ is the objective function defined in (21). To see this, first we recall from Proposition 6, the optimal transport map between two zero-mean Gaussians is $\mathbf{T}_{\mathbf{X} \to \mathbf{Y}} = \mathbf{M}(\mathbf{X}^{-1}\#(\mathbf{M}^{-1}\mathbf{Y}\mathbf{M}^{-1}))$ with $\mathbf{X}, \mathbf{Y}$ the respective covariance matrices. Now suppose $\mathbf{a} \in \mathbb{R}^n$ is a random Gaussian vector with mean zero and covariance $\mathbf{A}$ and define $\mathbf{x}_l = \mathbf{T}_{\mathbf{A} \to \mathbf{X}_l}\mathbf{a}$. From Proposition 6, $\mathbf{x}_l$ is Gaussian distributed with covariance $\mathbf{X}_l$ and

$$
F(\mathbf{A}) = \sum_{l=1}^{N} w_l\, d_{\mathrm{gbw}}^2(\mathbf{A}, \mathbf{X}_l) = \sum_{l=1}^{N} w_l\, \mathbb{E}\|\mathbf{a} - \mathbf{x}_l\|_{\mathbf{M}^{-1}}^2.
$$

In addition, we verify that $\mathbf{T}_{\mathbf{A}\to K(\mathbf{A})} = \sum_{l=1}^{N} w_l \mathbf{T}_{\mathbf{A}\to\mathbf{X}_l}$. That is,

$$
\begin{aligned}
\mathbf{T}_{\mathbf{A}\to K(\mathbf{A})} &= \mathbf{M}(\mathbf{A}^{-1}\#(\mathbf{M}^{-1}K(\mathbf{A})\mathbf{M}^{-1})) \\
&= \mathbf{M}\Big(\mathbf{A}^{-1}\#\big(\mathbf{A}^{-1/2}\big(\sum_{l=1}^{N} w_l\,(\mathbf{A}^{1/2}\mathbf{M}^{-1}\mathbf{X}_l\mathbf{M}^{-1}\mathbf{A}^{1/2})^{1/2}\big)^2 \mathbf{A}^{-1/2}\big)\Big) \\
&= \mathbf{M}\Big(\mathbf{A}^{-1/2}(\sum_{l=1}^{N} w_l(\mathbf{A}^{1/2}\mathbf{M}^{-1}\mathbf{X}_l\mathbf{M}^{-1}\mathbf{A}^{1/2})^{1/2})\mathbf{A}^{-1/2}\Big) \\
&= \sum_{l=1}^{N} w_l\,\mathbf{M}\Big(\mathbf{A}^{-1/2}(\mathbf{A}^{1/2}\mathbf{M}^{-1}\mathbf{X}_l\mathbf{M}^{-1}\mathbf{A}^{1/2})^{1/2}\mathbf{A}^{-1/2}\Big) = \sum_{l=1}^{N} w_l\mathbf{T}_{\mathbf{A}\to\mathbf{X}_l}.
\end{aligned}
$$

Denote $\bar{\mathbf{x}} := \sum_{l=1}^{N} w_l\mathbf{x}_l$. Then

$$
d_{\mathrm{gbw}}^2(\mathbf{A}, K(\mathbf{A})) = \mathbb{E}\|\mathbf{a}-\mathbf{T}_{\mathbf{A}\to K(\mathbf{A})}\,\mathbf{a}\|_{\mathbf{M}^{-1}}^2 = \mathbb{E}\|\mathbf{a}-\sum_{l=1}^{N} w_l\mathbf{x}_l\|_{\mathbf{M}^{-1}}^2 = \mathbb{E}\|\mathbf{a}-\bar{\mathbf{x}}\|_{\mathbf{M}^{-1}}^2.
$$

Notice $\bar{\mathbf{x}} = \mathbf{T}_{\mathbf{A}\to K(\mathbf{A})}\,\mathbf{a}$ is also a zero-mean Gaussian random vector with covariance $K(\mathbf{A})$. It follows that $d_{\mathrm{gbw}}^2(K(\mathbf{A}), \mathbf{X}_l) \le \mathbb{E}\|\bar{\mathbf{x}} - \mathbf{x}_l\|_{\mathbf{M}^{-1}}^2$. Next recall the variance formula for Euclidean random vector, i.e. $\mathrm{Var}(\mathbf{y}) = \mathbb{E}\|\mathbf{y} - \mathbb{E}[\mathbf{y}]\|^2 = \mathbb{E}\|\mathbf{y}\|^2 - \|\mathbb{E}[\mathbf{y}]\|^2 = \mathbb{E}\|\mathbf{x} - \mathbf{y}\|^2 - \|\mathbf{x} - \mathbb{E}[\mathbf{y}]\|^2$, for arbitrary $\mathbf{x}$. The analogue under Mahalanobis distance and finite average also holds, i.e., $\sum_{l=1}^{N} w_l\|\mathbf{x}_l - \bar{\mathbf{x}}\|_{\mathbf{M}^{-1}}^2 = \sum_{l=1}^{N} w_l\|\mathbf{a} - \mathbf{x}_l\|_{\mathbf{M}^{-1}}^2 - \|\mathbf{a} - \bar{\mathbf{x}}\|_{\mathbf{M}^{-1}}^2$. Finally, based on these results, we have

$$
\begin{aligned}
F(K(\mathbf{A})) = \sum_{l=1}^{N} w_l\,d_{\mathrm{gbw}}^2(K(\mathbf{A}), \mathbf{X}_l) &\le \sum_{l=1}^{N} w_l\,\mathbb{E}\|\bar{\mathbf{x}} - \mathbf{x}_l\|_{\mathbf{M}^{-1}}^2 \\
&= \sum_{l=1}^{N} w_l\,\mathbb{E}\|\mathbf{a} - \mathbf{x}_l\|_{\mathbf{M}^{-1}}^2 - \mathbb{E}\|\mathbf{a} - \bar{x}\|_{\mathbf{M}^{-1}}^2 \\
&\le F(\mathbf{A}) - d_{\mathrm{gbw}}^2(\mathbf{A}, K(\mathbf{A})).
\end{aligned}
$$

This suggests $F(K(\mathbf{A})) \le F(\mathbf{A})$ and hence together with the boundedness of $K(\mathbf{A})$, the sequence $\mathbf{A}_t$ converges. In the limit, we shall observe $F(K(\mathbf{A}_t)) = F(\mathbf{A}_t)$ when $t \to \infty$ and thus $d^2(\mathbf{A}, K(\mathbf{A})) = 0$. From the definition of $K(\mathbf{A})$ and the optimality condition, we conclude the limit point is $A(\mathbf{X}_{1:N}, \mathbf{w})$.

## C      Additional results and proofs for Section 3.2

### C.1     Geodesic convexity

Geodesic convexity is a generalization of standard convexity in the Euclidean space. It plays a crucial role in Riemannian optimization problems, where for geodesic convex problems, the convergence rates have been shown to be superior in many cases [52,57]. Consequently, geodesic convexity has been exploited to

develop better algorithms for machine learning applications such as Gaussian mixture models [26] and metric learning [56]. Below, we show some interesting classes of objective functions for SPD matrices that are geodesic convex under the GBW geometry.

A *geodesic convex set* $\mathcal{X} \subseteq \mathcal{M}$ requires, for any $x, y \in \mathcal{X}$, the distance minimizing geodesic $\gamma$ connecting the two points lie entirely in the set. A function $f : \mathcal{X} \to \mathbb{R}$ is called *geodesic convex* if, for any $x, y \in \mathcal{X}$, it satisfies that, for all $t \in [0, 1]$, $f(\gamma(t)) \leq (1 - t)f(x) + tf(y)$.

**Proposition 14.** *Suppose $\mathbf{A} \in \mathbb{S}_+^n$, the set of $n \times n$ semi-definite matrices, and let $\lambda^\downarrow : \mathbb{S}_{++}^n \to \mathbb{R}_+^n$ be the eigenvalue map that is decreasingly sorted and $h : \mathbb{R}_+ \to \mathbb{R}$ be a monotonically increasing and convex function. Then, the following functions $f_1(\mathbf{X}) = \mathrm{tr}(\mathbf{XA})$, $f_2(\mathbf{X}) = \mathrm{tr}(\mathbf{XAX})$, $f_3(\mathbf{X}) = -\log\det(\mathbf{X})$, $f_4(\mathbf{X}) = \sum_{j=1}^k h(\lambda_j^\downarrow(\mathbf{X}))$, $k \in [1, n]$, are geodesic convex under the GBW geometry for any choice of $\mathbf{M}$.*

### C.2   Proof of Proposition 5

Given a function $f : \mathcal{M} \to \mathbb{R}$, the Riemannian gradient at $x \in \mathcal{M}$, denoted by $\mathrm{grad}f(x)$, is the unique tangent vector satisfying $\langle \mathrm{grad}f(x), u \rangle_x = \mathrm{D}_u f(x)$, for any $u \in T_x\mathcal{M}$. $\mathrm{D}_u f(x)$ is the directional derivative. Riemannian Hessian at $x$, $\mathrm{Hess}f(x) : T_x\mathcal{M} \to T_x\mathcal{M}$ is defined as the Levi-Civita derivative of the Riemannian gradient, i.e., $\nabla \mathrm{grad}f(x)$.

*Proof (Proof of Proposition 5).* For the Riemannian gradient, we require

$$\mathrm{tr}(\nabla f(\mathbf{X})\mathbf{V}) = \frac{1}{2}\mathrm{tr}(\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathrm{grad}f(\mathbf{X})]\mathbf{V})$$

for any $\mathbf{V} \in T_{\mathbf{X}}\mathcal{M}_{\mathrm{gbw}}$. Thus, we have $\mathrm{grad}f(\mathbf{X}) = \mathcal{L}_{\mathbf{X},\mathbf{M}}^{-1}[2\nabla f(\mathbf{X})] = 2\mathbf{X}\nabla f(\mathbf{X})\mathbf{M} + 2\mathbf{M}\nabla f(\mathbf{X})\mathbf{X}$.

For the Riemannian Hessian, we have for any $\mathbf{U} \in T_{\mathbf{X}}\mathcal{M}$

$$\begin{aligned}
&\mathrm{Hess}f(\mathbf{X})[\mathbf{U}] = \nabla_{\mathbf{U}}\mathrm{grad}f(\mathbf{X}) \\
&= \mathrm{D}_{\mathbf{U}}\mathrm{grad}f(\mathbf{X}) - \{\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathrm{grad}f(\mathbf{X})]\mathbf{U}\}_{\mathrm{S}} - \{\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathrm{grad}f(\mathbf{X})\}_{\mathrm{S}} \\
&\quad + \{\mathbf{X}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathrm{grad}f(\mathbf{X})]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{M} + \mathbf{X}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathrm{grad}f(\mathbf{X})]\mathbf{M}\}_{\mathrm{S}} \\
&= \mathrm{D}_{\mathbf{U}}\mathrm{grad}f(\mathbf{X}) + \{4\mathbf{X}\{\nabla f(\mathbf{X})\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\}_{\mathrm{S}}\mathbf{M}\}_{\mathrm{S}} - \{2\mathbf{M}\nabla f(\mathbf{X})\mathbf{U}\}_{\mathrm{S}} \\
&\quad - \{\mathbf{M}\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{U}]\mathrm{grad}f(\mathbf{X})\}_{\mathrm{S}}, \quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad (19)
\end{aligned}$$

where we use $\mathcal{L}_{\mathbf{X},\mathbf{M}}[\mathbf{XUM} + \mathbf{MUX}] = \mathbf{U}$. Now we compute $\mathrm{D}_{\mathbf{U}}\mathrm{grad}f(\mathbf{X})$, which is

$$\begin{aligned}
\mathrm{D}_{\mathbf{U}}\mathrm{grad}f(\mathbf{X}) &= 2\mathrm{D}_{\mathbf{U}}(\mathbf{X}\nabla f(\mathbf{X})\mathbf{M} + \mathbf{M}\nabla f(\mathbf{X})\mathbf{X}) \\
&= 2\mathbf{U}\nabla f(\mathbf{X})\mathbf{M} + 2\mathbf{X}\nabla^2 f(\mathbf{X})[\mathbf{U}]\mathbf{M} + 2\mathbf{M}\nabla^2 f(\mathbf{X})[\mathbf{U}]\mathbf{X} + 2\mathbf{M}\nabla f(\mathbf{X})\mathbf{U} \\
&= 4\{\mathbf{M}\nabla f(\mathbf{X})\mathbf{U}\}_{\mathrm{S}} + 4\{\mathbf{M}\nabla^2 f(\mathbf{X})[\mathbf{U}]\mathbf{X}\}_{\mathrm{S}}. \quad\quad\quad (20)
\end{aligned}$$

Combining (20) with (19) completes the proof.

### C.3    Proof of Proposition 14

*Proof (Proof of Proposition 14).* To prove geodesic convexity for $f_1, f_2$, we require a second-order characterization of geodesic convexity. That is, a twice continuously differentiable function $f$ is geodesic convex if $\frac{d^2 f(\gamma(t))}{dt^2} \geq 0$ for all $t \in [0, 1]$. Now recall from Proposition 8 and the simplification in (8), the geodesic for GBW shares the same form as BW except for the value of polar factor $\mathbf{U}$. Nevertheless, the non-negativity of second-order derivatives does not depend on the choice of $\mathbf{U}$ according to the proof of Proposition 1 in [21]. Hence, we can follow the exact proof to show $f_1$ and $f_2$ are geodesic convex on the GBW geometry.

For $f_3$, we have

$$
\begin{aligned}
\log \det(\gamma(t)) &= 2 \log \det((1 - t)\mathbf{X}^{1/2} + t\mathbf{Y}^{1/2}\mathbf{U}) \\
&= 2 \log \det(((1 - t)\mathbf{M} + t\mathbf{Y}^{1/2}\mathbf{U}\mathbf{X}^{-1/2}\mathbf{M})\mathbf{M}^{-1}\mathbf{X}^{1/2}) \\
&\geq 2(1 - t) \log \det(\mathbf{M}) + 2t \log \det(\mathbf{Y}^{1/2}\mathbf{U}\mathbf{X}^{-1/2}\mathbf{M}) + 2 \log \det(\mathbf{M}^{-1}) \\
&\quad + 2 \log \det(\mathbf{X}^{1/2}) \\
&= 2t \log \det(\mathbf{Y}^{1/2}) - 2t \log \det(\mathbf{X}^{1/2}) + 2 \log \det(\mathbf{X}^{1/2}) \\
&= (1 - t) \log \det(\mathbf{X}) + t \log \det(\mathbf{Y}),
\end{aligned}
$$

where the inequality is due to the concavity of log-det on SPD matrices and from Lemma 2, we see $\mathbf{Y}^{1/2}\mathbf{U}\mathbf{X}^{-1/2}\mathbf{M} \succeq \mathbf{0}$.

Finally for $f_4$, the geodesic convexity simply follows from the result of $\mathbf{X} \star_t \mathbf{Y} \preceq (1 - t)\mathbf{X} + t\mathbf{Y}$ in Proposition 15 and Theorem 2.3 in [52]. $\qquad \square$

## D    Additional developments on the GBW geometry

### D.1    Results on geometric interpolation and barycenter

The geometric mean between symmetric positive definite matrices $\mathbf{X}$ and $\mathbf{Y}$ under the GBW geometry is the mid-point $\gamma(1/2)$ on the geodesic $\gamma$ that connects $\mathbf{X}$ to $\mathbf{Y}$. Following the notation in [8], we denote the interpolation of the generalized BW geodesic as $\mathbf{X} \star_t \mathbf{Y} := \gamma(t)$ derived in Proposition 8. We show an operator inequality between the interpolation on GBW and convex combination on the Euclidean space.

**Proposition 15.** *(Operator inequality) For any $\mathbf{X}, \mathbf{Y} \in \mathbb{S}_{++}^n$, we have $\mathbf{X} \star_t \mathbf{Y} \preceq (1 - t)\mathbf{X} + t\mathbf{Y}$, for $t \in [0, 1]$, where $\preceq$ denotes the Löwner partial order.*

An immediate result from this proposition is that $\log \det(\mathbf{X} \star_t \mathbf{Y}) \leq \log \det((1 - t)\mathbf{X} + t\mathbf{Y})$. This has implication in the application of Diffusion Tensor Imaging, where the larger determinant of interpolation of SPD matrices indicates the larger diffusion, known as the swelling effect, which is physically undesirable [4,46]. Because log det is geodesic concave on $\mathcal{M}_{\text{gbw}}$ (Proposition 14), the swelling effect still exists (unlike the affine-invariant or the log-Euclidean geometry), but the level of adverse effect is smaller compared to Euclidean metric.

Given a set of SPD matrices $\{\mathbf{X}_l\}_{l=1}^N$, the barycenter (or Riemannian center of mass) learning problem is

$$\min_{\mathbf{A} \in \mathbb{S}_{++}^n} F(\mathbf{A}) := \sum_{l=1}^N w_l d_{\mathrm{gbw}}^2(\mathbf{X}_l, \mathbf{A}), \tag{21}$$

with $\sum_{l=1}^N w_l = 1$. This is an extension of the Wasserstein barycenter of Gaussian measures [2,8]. Denote the minimizer as $A(\mathbf{X}_{1:N}, \mathbf{w}) := \arg\min_{\mathbf{A} \in \mathbb{S}_{++}^n} F(\mathbf{A})$. We can show, from matrix theory, that the minimizer is unique and is the solution to a specific nonlinear matrix equation. This generalizes the results in [8].

**Theorem 4 (Generalization of the result from [8]).** *The function $F(\mathbf{A})$ is strictly (Euclidean) convex in the convex cone of $\mathbb{S}_{++}^n$, which admits a unique GBW barycenter $A(\mathbf{X}_{1:N}, \mathbf{w})$. The barycenter is the solution to the equation $\mathbf{A}^{1/2}\mathbf{M}^{-1}\mathbf{A}^{1/2} = \sum_{l=1}^N w_l (\mathbf{A}^{1/2}\mathbf{M}^{-1}\mathbf{X}_l\mathbf{M}^{-1}\mathbf{A}^{1/2})^{1/2}$.*

Next, we show how to compute the barycenter by a fixed point iteration similar in [8,3]. Let

$$K(\mathbf{A}) := \mathbf{M}\mathbf{A}^{-1/2}\Big(\sum_{l=1}^N w_l (\mathbf{A}^{1/2}\mathbf{M}^{-1}\mathbf{X}_l\mathbf{M}^{-1}\mathbf{A}^{1/2})^{1/2}\Big)^2 \mathbf{A}^{-1/2}\mathbf{M},$$

and perform the iteration update by $\mathbf{A}_{t+1} = K(\mathbf{A}_t)$. We can show this update converges to $A(\mathbf{X}_{1:N}, \mathbf{w})$, formalized in the following Theorem.

**Theorem 5.** *Initialize $\mathbf{A}_0 \in \mathbb{S}_{++}^n$ randomly and consider the update $\mathbf{A}_{t+1} = K(\mathbf{A}_t)$. Then $\lim_{t \to \infty} \mathbf{A}_t = A(\mathbf{X}_{1:N}, \mathbf{w})$.*

### D.2    Robust GBW distance

In this section, we show that the connection of the GBW distance with a class of projection robust Wasserstein distances between zero-centered Gaussians. This may be of independent interest.

Robust Wasserstein distances [45,28] may help mitigate the sample complexity of Wasserstein distances, which may grow exponentially in dimension [16,17,55]. Given two $n$-dimensional measures $\mu, \nu$, the projection robust Wasserstein distance [45,28] is computed as follows:

$$\mathcal{P}_d(\mu, \nu) = \sup_{\mathbf{W}:\mathbf{W}^\top\mathbf{W}=\mathbf{I}} \inf_{\gamma \sim \Gamma(\mu,\nu)} \int \|\mathbf{W}^\top(\mathbf{x} - \mathbf{y})\|^2 d\gamma(\mathbf{x}, \mathbf{y}),$$

where $\mathbf{W} \in \mathbb{R}^{n \times d}$ $(d \leq n)$ is a projection matrix which is learned over the given samples. When $\mu$ and $\nu$ are zero-centered Gaussians with covariance matrices $\mathbf{X}$ and $\mathbf{Y}$, respectively, this reduces to

$$\mathcal{P}_d(\mu = \mathcal{N}(\mathbf{0}, \mathbf{X}), \nu = \mathcal{N}(\mathbf{0}, \mathbf{Y})) = \tag{22}$$
$$\sup_{\mathbf{W}:\mathbf{W}^\top\mathbf{W}=\mathbf{I}} \mathrm{tr}(\mathbf{W}\mathbf{W}^\top\mathbf{X}) + \mathrm{tr}(\mathbf{W}\mathbf{W}^\top\mathbf{Y}) - 2\mathrm{tr}(\mathbf{X}^{1/2}\mathbf{W}\mathbf{W}^\top\mathbf{Y}\mathbf{W}\mathbf{W}^\top\mathbf{X}^{1/2})^{1/2}$$

based on Proposition 2. If $\mathbf{W}^*$ is an optimal solution of (22), we also have the following equivalence: $\mathcal{P}_d(\mu = \mathcal{N}(\mathbf{0}, \mathbf{X}), \nu = \mathcal{N}(\mathbf{0}, \mathbf{Y})) = d_{\text{gbw}}^2(\mathbf{X}, \mathbf{Y})$ for $\mathbf{M}^{-1} = \mathbf{W}^*(\mathbf{W}^*)^\top$. Hence, for a specific choice of $\mathbf{M}^{-1}$, the GBW distance may be interepreted as a projection robust Wasserstein distance between zero-centered Gaussians.

Based on the above discussion, we now define a class of robust Wasserstein distances $d_{\text{rgbw}}$ for $\mathbf{M}^{-1} \succ \mathbf{0}$ as

$$d_{\text{rgbw}}^2(\mathbf{X}, \mathbf{Y}) = \max_{\mathbf{M}^{-1} \in \mathcal{C}} d_{\text{gbw}}^2(\mathbf{X}, \mathbf{Y}) = \max_{\mathbf{S} \in \mathcal{C}} \text{tr}(\mathbf{SX}) + \text{tr}(\mathbf{SY}) - 2\text{tr}(\mathbf{X}^{1/2}\mathbf{SYSX}^{1/2})^{1/2}$$
(23)

for a closed convex set $\mathcal{C} \subseteq \mathbb{S}_{++}^n$. We emphasize the maximization of $\mathbf{S}$ over the set $\mathcal{C}$. Below we show that (23) is a distance metric.

**Proposition 16.** *The robust GBW distance (23) in the set $\mathcal{C} \subseteq \mathbb{S}_{++}^n$ is a distance metric.*

*Proof (Proof of Proposition 16).* From (23), we see $d_{\text{rgbw}}^2(\mathbf{X}, \mathbf{Y}) \geq 0$ and is clearly symmetric. The triangle inequality also easily follows as shown below. Let

$$\mathbf{S}^* = \arg\max_{\mathbf{S} \in \mathcal{C}} d_{\text{gbw}}^2(\mathbf{X}, \mathbf{Y}).$$
(24)

Therefore, from (24), we have

$$\begin{aligned}
d_{\text{rgbw}}(\mathbf{X}, \mathbf{Y}) &= d_{gbw}(\mathbf{X}, \mathbf{Y}) \text{ for } \mathbf{S}^* \\
&\leq d_{gbw}(\mathbf{X}, \mathbf{Z}) + d_{gbw}(\mathbf{Z}, \mathbf{Y}) \text{ for } \mathbf{S}^* \text{ as GBW is a distance} \\
&\leq (\max_{\mathbf{S}_1 \in \mathcal{C}} d_{gbw}(\mathbf{X}, \mathbf{Z}) \text{ for } \mathbf{S}_1) + (\max_{\mathbf{S}_2 \in \mathcal{C}} d_{gbw}(\mathbf{Z}, \mathbf{Y}) \text{ for } \mathbf{S}_2) \\
&= d_{\text{rgbw}}(\mathbf{X}, \mathbf{Z}) + d_{\text{rgbw}}(\mathbf{Z}, \mathbf{Y}),
\end{aligned}$$

where $\mathbf{X}$, $\mathbf{Y}$, and $\mathbf{Z}$ are SPD matrices. Finally, the identity of indiscernibles property is satisfied as the robust GBW distance is based on the GBW distance (which itself satisfies the property). This completes the proof.

### D.3    GBW geometry and metric learning

The problem of metric learning amounts to learning a suitable Mahalanobis (symmetric positive, and possibly, semi-definite) matrix from pairs of similarity and dissimilarity information, e.g., in a classification task [56,22,20,24,30,31].

A particular formulation of interest is based on the objective function proposed in [22]. Specifically, given a set of data-target pairs $\{\mathbf{X}_i, t_i\}, \mathbf{X}_i \in \mathbb{S}_{++}^n$ and $t_i$ categorical, we define the class adjacency matrix $\mathbf{A}_{ij} = 1$ if sample $i, j$ are from the same class (i.e., $t_i = t_j$) and $\mathbf{A}_{ij} = -1$ otherwise. To this end, the objective function is given as

$$\min_{\mathbf{S} \succeq \mathbf{0}} \sum_{i,j}^N \log(1 + \exp(\mathbf{A}_{ij}(\text{tr}(\mathbf{SX}_i) + \text{tr}(\mathbf{SX}_j) - 2\text{tr}(\mathbf{X}_i^{1/2}\mathbf{SX}_j\mathbf{SX}_i^{1/2})^{1/2}))). \quad (25)$$

It should be emphasized that the objective function in (25) is formulated by directly making use of the BW distance in the objective function of [22]. However, from the definition of the GBW distance (4) between $\mathbf{X}_i$ and $\mathbf{X}_j$ and by taking $\mathbf{M}^{-1} = \mathbf{S}$, we observe that the problem (25) may be equivalently rewritten as

$$\min_{\mathbf{S} \succeq \mathbf{0}} \sum_{i,j}^{N} \log(1 + \exp(\mathbf{A}_{ij} d_{\text{gbw}}^2(\mathbf{X}_i, \mathbf{X}_j))).$$

This suggests that the GBW geometry naturally captures the metric learning properties of the space. Note that $\mathbf{S}$ can be arbitrary semi-definite matrix and one usually parameterizes $\mathbf{S} = \mathbf{W}\mathbf{W}^\top$, where $\mathbf{W}$ is a matrix of size $n \times d$. Similar to Section E, we usually consider $d \ll n$ for practical considerations.

## E    Additional experiments on geometry-aware principal component analysis (PCA)

In this section, we explore the connection of the GBW distance and geometry-aware principal component analysis.

***Problem formulation:*** Geometry-aware principal component analysis (PCA) for SPD matrices extends the classical PCA to manifolds by maximizing the deviation from the reduced SPD matrices to the reduced barycenter [25,24,31]. Using the BW distance, the PCA objective is formulated naturally as the GBW distance between matrices, where $\mathbf{M}^{-1}$ is parameterized as $\mathbf{W}\mathbf{W}^\top$ with $\mathbf{W} \in \mathbb{R}^{n \times d}$. Note that $\mathbf{M}^{-1}$ is low rank, therefore, does not strictly fall under the generalized metric. Nevertheless, we can make use of the GBW distance expression and substitute low-rank paramterized $\mathbf{M}^{-1}$.

Consequently, the objective function is

$$\max_{\mathbf{M}^{-1} = \mathbf{W}\mathbf{W}^\top : \mathbf{W}^\top \mathbf{W} = \mathbf{I}} \sum_{i=1}^{N} d_{\text{gbw}}^2(\mathbf{X}_i, \bar{\mathbf{X}}) = \max_{\mathbf{W} : \mathbf{W}^\top \mathbf{W} = \mathbf{I}} \sum_{i=1}^{N} d_{\text{bw}}^2(\mathbf{W}^\top \mathbf{X}_i \mathbf{W}, \mathbf{W}^\top \bar{\mathbf{X}} \mathbf{W})$$

for samples $\mathbf{X}_i \in \mathbb{S}_{++}^n, i = 1, \ldots, N$, where $\bar{\mathbf{X}} = \arg\min \sum_{i=1}^{N} d_{\text{bw}}^2(\mathbf{X}_i, \mathbf{C})$ is the barycenter in the original space. The constraint of column orthonormality on $\mathbf{W}$, i.e., $\mathbf{W}^\top \mathbf{W} = \mathbf{I}$, ensures that $\mathbf{W}$ projects the covariance matrices onto a $d$-dimensional space. In many practical scenarios, $d$ is often chosen to be much less than $n$, i.e., $d \ll n$.

***Tasks:*** For the application of geometry-aware PCA, we consider two vision tasks, i.e., image set classification and video-based face recognition. Following the pre-processing steps in [24,31], we treat each vectorized image (or video frame) as a sample in the set and compute the sample covariance to represent the entire image set (or a video). The task is to classify each image set or video represented by a covariance SPD matrix.

**Table 3.** Summary statistics for `MNIST`, `ETH`, `YTC` datasets

|        | SPD samples | SPD Dim | # Class |
|--------|-------------|---------|---------|
| `MNIST` | 835         | 100     | 10      |
| `ETH`   | 80          | 100     | 8       |
| `YTC`   | 194         | 100     | 9       |

**Table 4.** Geometry-aware PCA average classification accuracy (%). GBW allows lower dimensional projection with accuracy comparable to that in the original dimension.

|        | AI    | LE    | BW    | GBW | | | | | |
|--------|-------|-------|-------|-------|--------|--------|--------|--------|--------|
|        |       |       |       | $d=5$ | $d=10$ | $d=30$ | $d=50$ | $d=70$ | $d=90$ |
| `MNIST` | 100   | 100   | 100   | 99.33 | 100    | 100    | 100    | 100    | 100    |
| `ETH`   | 76.25 | 84.50 | 87.75 | 80.75 | 84.75  | 86.75  | 88.00  | 87.75  | 87.75  |
| `YTC`   | 74.70 | 79.00 | 76.40 | 60.60 | 72.40  | 76.50  | 76.00  | 76.30  | 76.40  |

***Datasets:*** Three real-world datasets are considered, including the MNIST handwritten digits (`MNIST`) [34], ETH-80 object (`ETH`) [37], and YouTube Celebrities (`YTC`) [32] datasets. To process `MNIST` dataset, we use $42\,000$ training samples, and, for each class, we partition the samples into subgroups randomly, each containing 50 images. Then for each subgroup, the covariance matrix is computed. `ETH` dataset contains image sets of 8 objects, each with 10 subclasses. The 80 subgroups are processed accordingly. `YTC` is a collection of low-resolution videos of celebrities. Due to the sparsity of the dataset, we only consider 9 persons with video number greater than 15. All images or video frames are resized to $10 \times 10$ and the SPD matrix generated as the covariance is of size $100 \times 100$. The statistics of all the considered datasets are in Table 3.

***Experimental setup:*** As discussed in the above problem formulation, our aim is to find the transformation matrix $\mathbf{W} \in \mathbb{R}^{n \times d}$. To validate the effectiveness of dimensionality reduction under the GBW geometry, we perform nearest neighbour classification on the reduced data matrix $\mathbf{W}^\top \mathbf{X}_i \mathbf{W}, i = 1, \ldots, N$. The reduced dimension $d$ is a hyperparameter, and we, therefore, present classification accuracy with $d = \{5, 10, 30, 50, 70, 90\}$. Given that the sample size may be small for some classes, for each class, we take 50% as the training set and the rest as the test set. Such a random splitting is repeated ten times and we report the average accuracy in Table 4, where we also report results with the affine-invariant (AI) and Log-Euclidean (LE) [4] distances as benchmarks. We use the Riemannian trust region method to solve the maximization problem in Section E.

***Results:*** In Table 4, we observe that the classification performance under various choices of $d$ for the GBW distance does not largely degrade, which

suggests the global properties of SPD samples can be well-preserved even with a lower-dimensional representation. This also suggests that GBW is a better modeling approach than BW for the geometric PCA problem.