

---

# Riemannian Diffusion Adaptation for Distributed Optimization on Manifolds

---

Xiuheng Wang<sup>1</sup> Ricardo Borsoi<sup>1</sup> Cédric Richard<sup>2</sup> Ali H. Sayed<sup>3</sup>

## Abstract

Online distributed optimization is particularly useful for solving optimization problems with streaming data collected by multiple agents over a network. When the solutions lie on a Riemannian manifold, such problems become challenging to solve, particularly when efficiency and continuous adaptation are required. This work tackles these challenges and devises a diffusion adaptation strategy for decentralized optimization over general manifolds. A theoretical analysis shows that the proposed algorithm is able to approach network agreement after sufficient iterations, which allows a non-asymptotic convergence result to be derived. We apply the algorithm to the online decentralized principal component analysis problem and Gaussian mixture model inference. Experimental results with both synthetic and real data illustrate its performance.

## 1. Introduction

In the decentralized setting, this work deals with the multi-agent optimization problem seeking *consensus* on a Riemannian manifold  $\mathcal{M}$ :

$$\min_{w \in \mathcal{M}} \frac{1}{K} \sum_{k=1}^K J_k(w), \quad (1)$$

where  $J_k : \mathcal{M} \rightarrow \mathbb{R}$  is a local risk function defined for each agent by  $J_k(w) = \mathbb{E}_{\mathbf{x}_k} \{Q(w; \mathbf{x}_k)\}$  in terms of the expectation of some loss function  $Q(w; \mathbf{x}_k)$ . The computation of  $J_k(w)$  is over the unknown distribution of the data  $\mathbf{x}_k$ , which makes it necessary to use a stochastic approximation for the gradient vector based on a set of independent realizations  $\mathbf{x}_{k,t}$ , observed sequentially over time. A wide range

of applications in machine learning, signal processing, and control can be written in the form of (1), including principal component analysis (PCA) (Cunningham & Ghahramani, 2015; Zhang et al., 2016), parameter estimation for Gaussian mixture models (GMM) (Hosseini & Sra, 2015; Collas et al., 2023), low-rank matrix completion (Boumal & Absil, 2011; Vandereycken, 2013), and deep neural networks with orthogonal constraints (Vorontsov et al., 2017).

Decentralized optimization in Euclidean spaces has been extensively studied. As such, one may consider converting the constraint  $w \in \mathcal{M}$  into a cost function in the Euclidean space and solve (1) in a composite setting. Then, a distributed proximal gradient-type algorithm can be applied if the projection onto the manifold is available. However, previous studies, e.g., (Bianchi & Jakubowicz, 2012; Di Lorenzo & Scutari, 2016; Zeng & Yin, 2018), require a convex regularizer or at least the convexity of its domain. In addition, for certain manifolds, the dimension of an embedded Euclidean space can be relatively high according to the Whitney Embedding Theorem (Lee, 2013). Due to the non-convexity and non-linearity of certain  $\mathcal{M}$ , these decentralized algorithms may fail when dealing with problem (1).

In response to these challenges, this paper aims to introduce a general framework to solve (1) by developing fully decentralized optimization on manifolds, which directly operates on  $\mathcal{M}$  by exploiting its inherent geometry. Our main contributions are as follows:

- 1. Riemannian diffusion adaptation:** We devise a Riemannian diffusion adaptation strategy, which is fully intrinsic and thus can be applied to general manifolds<sup>1</sup>. It comprises a sequence of efficient *adaptation* and *combination* steps. In the adaptation step, a Riemannian stochastic gradient descent (R-SGD) method is used to estimate the local solution at each agent. In the combination step, the local estimates of the neighboring agents are combined on the tangent space of the manifold.
- 2. Theoretical analysis:** We provide a theoretical analysis for the performance of the proposed Riemannian diffusion adaptation strategy with a constant step size.

---

Part of this work was done while Xiuheng Wang was a PhD student at Université Côte d’Azur. <sup>1</sup>Université de Lorraine, CNRS, CRAN, France <sup>2</sup>Université Côte d’Azur, CNRS, OCA, France <sup>3</sup>École Polytechnique Fédérale de Lausanne, Switzerland. Correspondence to: Xiuheng Wang <dr.xiuheng.wang@gmail.com>.

---

<sup>1</sup>This strategy is not designed for a specific manifold as it does not require an embedding of  $\mathcal{M}$  into a Euclidean space.

We establish that all agents will approximately converge to a network agreement (Theorem 5.12) in the sense of a decreasing geodesic distance between their estimates over the iterations. Additionally, we establish a curvature-dependent and non-asymptotic convergence result with a proper design of a Lyapunov function (Theorem 5.15).

3. **Application to various manifolds:** We tailor our algorithm to suit two instances of Riemannian manifolds, i.e., the Grassmann manifold and a product manifold involving the manifold of symmetric positive definite (SPD) matrices. We apply our algorithm to the online distributed PCA and parameter estimation of GMMs through numerical experiments on synthetic and real-world datasets. Experimental results show its performance compared to non-cooperative, consensus, and centralized solutions.

## 2. Related work

In this section, we review related works on decentralized optimization in Euclidean spaces and the recent advances in (decentralized) optimization on Riemannian manifolds.

**Decentralized optimization in Euclidean spaces:** Distributed optimization in Euclidean spaces has been extensively studied, including incremental (Blatt et al., 2007), consensus (Nedic et al., 2010) and diffusion (Chen & Sayed, 2012; Sayed et al., 2013) strategies. In particular, diffusion strategies have been demonstrated in (Sayed, 2014; Chen & Sayed, 2015) to offer improved performance and stability guarantees under constant step-size learning and adaptive scenarios. Recently, decentralized optimization has been extensively studied in non-convex environments (Bianchi & Jakubowicz, 2012; Di Lorenzo & Scutari, 2016; Lian et al., 2017; Tatarenko & Touri, 2017; Zeng & Yin, 2018; Wang et al., 2019; Vlaski & Sayed, 2021).

**Optimization on manifolds:** Riemannian optimization has garnered significant interest as it considers the geometry of manifolds, recently presented in detail in the books (Absil et al., 2009) and (Boumal, 2023). Of particular interest are stochastic optimization methods due to their efficiency and scalability. The first asymptotic convergence analysis of R-SGD was provided in (Bonnabel, 2013), highlighting diverse applications such as PCA. The first global convergence results for first-order Riemannian optimization with geodesic convexity were obtained in (Zhang & Sra, 2016). The finite-sum, stochastic setting has been further investigated in (Zhang et al., 2016; Sato et al., 2019) for variance reduction. The work (Tripuraneni et al., 2018) constructed and analyzed a variant of R-SGD that generalizes the classical Polyak-Ruppert iterate-averaging scheme. Several

stochastic Riemannian Frank-Wolfe methods were introduced in (Weber & Sra, 2022). More recently, the behavior of various stochastic optimization algorithms around saddle points in geodesically non-convex functions was studied in (Hsieh et al., 2024). In (Wang et al., 2024a), the non-asymptotic convergence of R-SGD with constant step sizes was studied and applied to the change point detection on manifolds.

**Decentralized optimization on manifolds:** The literature on decentralized optimization on Riemannian manifolds can be roughly divided into *extrinsic* and *intrinsic* methods.

The extrinsic methods are based on *induced arithmetic mean* (Sarlette & Sepulchre, 2009), and rely on the specific embedding of the manifolds in Euclidean spaces (where traditional Euclidean consensus can be employed), which is often studied for specific manifolds. For stochastic optimization, the incremental and consensus strategies have been extended to decentralized R-SGD-type on the unit sphere (Wang et al., 2023) and Stiefel manifolds (Chen et al., 2021), respectively. For the deterministic case, an augmented Lagrangian method (Wang & Liu, 2022) and a type of conjugate gradient method (Chen et al., 2024) were also designed for decentralized optimization on the Stiefel manifold. In addition, a consensus strategy has also been extended to compact submanifolds (Deng & Hu, 2023).

The intrinsic methods are based on *Fréchet mean* (Tron et al., 2012) (or center of mass) and developed with the inherent geometry of manifolds, such as geodesic distance, Riemannian gradient, and exponential mapping. These methods can be studied on more general manifolds including, but not limited to, the unit sphere, Stiefel manifolds, Grassmann manifolds, and the manifold of SPD matrices. Different distributed strategies (Tron et al., 2012; Kraisler et al., 2023a;b) were studied to achieve network agreement on manifolds. Another distributed strategy (Shah, 2017) was considered to solve (1) with a diminishing step size and two rounds of communication in each iteration. However, few methods have investigated the diffusion strategy on manifolds, though it has been proven to have superior properties in Euclidean spaces, especially in continuous learning and adaptive scenarios. A work extending the diffusion strategy to manifolds was introduced in (Wang et al., 2024b), but the algorithm is inefficient due to inner-loop optimization<sup>2</sup> and does not have any theoretical analyses.

Recently, another branch of distributed optimization on manifolds considering a central server was also investigated, with settings of communication efficiency (Huang & Pan, 2020) and federated learning (Li & Ma, 2023; Huang et al., 2024a;b).

<sup>2</sup>We further support this claim with a numerical evaluation in Appendix D.1

### 3. Background

This section introduces some basic concepts of Riemannian geometry, focusing on the essential tools for optimization on manifolds. Detailed presentations can be found in (Absil et al., 2009) and (Boumal, 2023).

A *Riemannian manifold*  $(\mathcal{M}, g)$  is a constrained set  $\mathcal{M}$  endowed with a *Riemannian metric*  $g_x(\cdot, \cdot) : T_x\mathcal{M} \times T_x\mathcal{M} \rightarrow \mathbb{R}$ , defined for every point  $x \in \mathcal{M}$ , with  $T_x\mathcal{M}$  the so-called *tangent space* of  $\mathcal{M}$  at  $x$ . A *geodesic*  $\gamma_v : [0, 1] \rightarrow \mathcal{M}$  is the curve of minimal length linking two points  $x, y \in \mathcal{M}$  such that  $x = \gamma_v(0)$  and  $y = \gamma_v(1)$ , with  $v \in T_x\mathcal{M}$  the velocity of  $\gamma_v$  at 0 denoted by  $\dot{\gamma}_v(0)$ . The *geodesic distance*  $d(\cdot, \cdot) : \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}$  is defined as the length of the geodesic linking two points  $x, y \in \mathcal{M}$ . It satisfies all the conditions to be a metric.

The *exponential map*  $w = \exp_x(v)$  is defined as the point  $w \in \mathcal{M}$  located on the unique geodesic  $\gamma_v(t)$  with endpoints  $x = \gamma_v(0)$ ,  $w = \gamma_v(1)$  and velocity  $v = \dot{\gamma}_v(0)$ . Consider a smooth function  $f : \mathcal{M} \rightarrow \mathbb{R}$ . The *Riemannian gradient* of  $f$  at  $x \in \mathcal{M}$  is defined as the unique tangent vector  $\nabla f(x) \in T_x\mathcal{M}$  satisfying  $\frac{d}{dt}\big|_{t=0} f(\exp_x(tv)) = \langle \nabla f(x), v \rangle_x$ , for all  $v \in T_x\mathcal{M}$ . The *Riemannian Hessian* of  $f$  at  $x$  is an operator  $\nabla_x^2 f$  such that  $\frac{d}{dt}\big|_{t=0} \langle \nabla f(\exp_x(tv)), \nabla f(\exp_x(tv)) \rangle_x = 2\langle \nabla f(x), (\nabla_x^2 f)v \rangle_x$ .

### 4. Algorithm development

Let us define the product manifold  $\mathcal{M}^K \triangleq \mathcal{M} \times \cdots \times \mathcal{M}$ , which is the  $K$ -fold Cartesian product of  $\mathcal{M}$  with itself. We also define  $\mathbf{w} \triangleq \text{col}\{\mathbf{w}_1, \dots, \mathbf{w}_K\}$  to indicate a point on  $\mathcal{M}^K$ . The decentralized optimization problem (1) contains an implicit consensus: the individual models are required to be common on the manifolds, i.e.,  $\mathbf{w}_k = w, \forall k$ . One can encourage consensus on manifolds by penalizing pairwise differences between connected agents. Let us represent the  $K$  agents as the nodes of a graph  $\mathcal{G}$ . With a natural generalization of the Euclidean case, we consider the geodesic distance-based consensus problem (Tron et al., 2012), i.e., minimization of the penalty  $P(\mathbf{w}) \triangleq \sum_{k=1}^K P_k(\mathbf{w}_k)$  where  $P_k(\mathbf{w}_k) \triangleq \frac{1}{2} \sum_{\ell=1}^K c_{\ell k} d^2(\mathbf{w}_k, \mathbf{w}_\ell)$  and  $c_{\ell k} \triangleq [C]_{\ell k}$ , with  $C$  a weighted adjacency matrix of the graph  $\mathcal{G}$ , representing the strength of link between each pair of agents. For a connected graph,  $P(\mathbf{w}) = 0$  if and only if  $\{\mathbf{w}_k\}_{k=1}^K$  are equal for all  $k$ . This results in the following optimization problem with a constraint:

$$\min_{\mathbf{w} \in \mathcal{M}^K} J(\mathbf{w}) \quad \text{s.t.} \quad P(\mathbf{w}) = 0, \quad (2)$$

where  $J(\mathbf{w}) \triangleq \frac{1}{K} \sum_{k=1}^K J_k(\mathbf{w}_k)$ . To minimize the global cost function in (2), we follow the diffusion adaptation strategy in Euclidean spaces (Sayed et al., 2014; Yuan et al., 2018; Vlaski et al., 2023), and appeal to an incremental gradient descent argument. We first apply an R-SGD to

---

#### Algorithm 1 Riemannian Diffusion Adaptation

---

**Input:** Step sizes  $\mu, \alpha$ , graph adjacency matrix  $C$ .  
 Initialize  $\{\mathbf{w}_{k,0}\}$  for all  $k$  with a random point on  $\mathcal{M}$ .  
**for**  $t = 1, 2, \dots$  **do**  
   **for** each agent  $k$  **do**  
      $\phi_{k,t} = \exp_{\mathbf{w}_{k,t-1}}(-\mu \widehat{\nabla J}_k(\mathbf{w}_{k,t-1}))$ ;  
      $\mathbf{w}_{k,t} = \exp_{\phi_{k,t}}(\alpha \sum_{\ell=1}^K c_{\ell k} \exp_{\phi_{k,t}}^{-1}(\phi_{\ell,t}))$ ;  
   **end for**  
**end for**

---

the risk  $J(\mathbf{w})$  and subsequently descend along the penalty  $P(\mathbf{w})$ . Using node-level quantities we have

$$\phi_{k,t} = \exp_{\mathbf{w}_{k,t-1}}(-\mu \widehat{\nabla J}_k(\mathbf{w}_{k,t-1})), \quad (3)$$

$$\mathbf{w}_{k,t} = \exp_{\phi_{k,t}}(-\alpha \nabla P_k(\phi_{k,t})), \quad (4)$$

where  $\mu$  and  $\alpha$  are step sizes,  $\widehat{\nabla J}_k$  is a stochastic approximation of the Riemannian gradient of  $J_k$ , and  $\nabla P_k$  can be computed in an explicit form (Afsari et al., 2013), given by

$$\begin{aligned} \nabla P_k(\phi_{k,t}) &= \frac{1}{2} \sum_{\ell=1}^K c_{\ell k} \nabla d^2(\phi_{k,t}, \phi_{\ell,t}) \\ &= - \sum_{\ell=1}^K c_{\ell k} \exp_{\phi_{k,t}}^{-1}(\phi_{\ell,t}). \end{aligned} \quad (5)$$

The Riemannian diffusion adaptation strategy, summarized in Algorithm 1, contains two steps: an adaptation step (3) where agent  $k$  uses its own data  $\mathbf{x}_{k,t-1}$  to update its solution  $\phi_{k,t}$  and a combination step (4) where the intermediate estimates  $\{\phi_{\ell,t}\}$  are combined, on the tangent space of  $\phi_{k,t}$ , according to the weighting coefficients  $\{c_{\ell k}\}$  in (5) to obtain the estimate  $\mathbf{w}_{k,t}$ . Note that in the special case that  $\mathcal{M}$  is a Euclidean space, we can take  $\exp_x(v)$  as vector addition of  $x + v$ , and our algorithm reduces to the diffusion adaptation algorithm in the Euclidean space (Chen & Sayed, 2012; Sayed et al., 2013). Here we emphasize that the local update of each agent in (3) is performed by stochastic Riemannian optimization with constant step size, which plays an important role in tasks in need of continuous learning and adaptation (Sayed et al., 2013; Sayed, 2014).

### 5. Theoretical analysis

In this section, we analyze the convergence of Algorithm 1 in the constant step size setting.

In analyzing the dynamics of the distributed algorithm (3) and (4), it is useful to introduce the following stacked vector notation by collecting variables from across the network:

$$\begin{aligned} \mathbf{w}_t &\triangleq \text{col}\{\mathbf{w}_{1,t}, \dots, \mathbf{w}_{K,t}\} \in \mathcal{M}^K \\ \widehat{\nabla J}(\mathbf{w}_t) &\triangleq \text{col}\{\widehat{\nabla J}_1(\mathbf{w}_{1,t}), \dots, \widehat{\nabla J}_K(\mathbf{w}_{K,t})\} \in T_{\mathbf{w}_t} \mathcal{M}^K \end{aligned}$$

$$\begin{aligned}\phi_t &\triangleq \text{col}\{\phi_{1,t}, \dots, \phi_{K,t}\} \in \mathcal{M}^K \\ \nabla P(\phi_t) &\triangleq \text{col}\left\{-\sum_{\ell=1}^K c_{\ell 1} \exp_{\phi_{1,t}}^{-1}(\phi_{\ell,t}), \dots, \right. \\ &\quad \left.-\sum_{\ell=1}^K c_{\ell K} \exp_{\phi_{K,t}}^{-1}(\phi_{\ell,t})\right\} \in T_{\phi_t} \mathcal{M}^K\end{aligned}$$

where  $\text{col}\{\cdot\}$  is obtained by stacking the arguments columnwise and  $T_x \mathcal{M}^K$  is the tangent space of  $\mathcal{M}^K$  at  $x$ , see Proposition 3.20 in (Boumal, 2023). We can then write (3) and (4) compactly as

$$\phi_t = \exp_{w_{t-1}}(-\mu \widehat{\nabla J}(w_{t-1})), \quad (6)$$

$$w_t = \exp_{\phi_t}(-\alpha \nabla P(\phi_t)). \quad (7)$$

Step (7) can be regarded as a one-step Riemannian gradient descent with a step size  $\alpha$  to approximate a global minimum of  $P(\phi)$ , belonging to the *consensus submanifold*  $\mathcal{A}$ , defined as

$$\mathcal{A} \triangleq \{\phi \in \mathcal{M}^K \mid \phi_i = \phi_j, \forall i, j\}. \quad (8)$$

We start by introducing some technical assumptions and existing auxiliary results before presenting new results.

### 5.1. Assumptions and auxiliary results

Let us denote the *convexity submanifold* (Tron et al., 2012) of product manifolds  $\mathcal{M}^K$  as  $\mathcal{B} \subseteq \mathcal{M}^K$ . We introduce the following standard assumptions in the literature on Riemannian optimization:

**Assumption 5.1 (Regularization on manifold).** (Bonnabel, 2013; Zhang et al., 2016; Tripuraneni et al., 2018; Afsari, 2011) (a) The sequences  $\{\phi_t\}_{t \geq 0}$  and  $\{w_t\}_{t \geq 0}$  generated by the algorithm stay continuously in  $\mathcal{B}$ , and  $J$  attains its optimum  $w^*$  in  $\mathcal{B}$ ; (b) the sectional curvature in  $\mathcal{B}$  is *upper* bounded by  $\kappa_{\max}$ ; (c) the sectional curvature in  $\mathcal{B}$  is *lower* bounded by  $\kappa_{\min}$ ; and (d)  $\mathcal{B}$  is compact, and the diameter of  $\mathcal{B}$  is bounded by  $D$ , that is,  $\max_{x,y \in \mathcal{B}} d(x,y) \leq D$ ; (e)  $D < D^*$ , where  $D^*$  is defined as  $D^* \triangleq \min(\text{inj}(\mathcal{M}), \frac{\pi}{\sqrt{\kappa_{\max}}})$  with  $\text{inj}(\mathcal{M})$  is the injectivity radius of  $\mathcal{M}$ , which implies that the exponential map is invertible within  $\mathcal{B}$ .

Also, it is necessary to assume some properties of the weighted adjacency matrix  $C$  according to which the agents interact over the graph topology  $\mathcal{G}$ . In addition to the direct assumptions on  $\mathcal{G}$  (e.g., left-stochastic) in Euclidean space (Chen & Sayed, 2012; Sayed et al., 2013), for distributed optimization on manifolds, we also make the following assumptions on the eigenvalues of the Riemannian Hessian of  $P$ , whose computation involves  $C$ , see Subsection 2.1.3 in (Afsari et al., 2013) and Proposition 8 in (Tron et al., 2012) for examples.

**Assumption 5.2 (Regularization on graph).** (Chen & Sayed, 2012; Sayed et al., 2013; Afsari et al., 2013) Assume that the undirected  $\mathcal{G}$  is connected and its adjacency matrix  $C$  is left-stochastic, i.e.,  $c_{\ell k} \geq 0$ ,  $\sum_{\ell=1}^K c_{\ell k} = 1$  for each agent  $k$ , denote lower and upper bounds on the eigenvalues of the Hessian of  $P$  in  $\mathcal{B}$  as  $h_{\min}$  and  $h_{\max}$ , and suppose that  $h_{\min} \geq 0$ .

Note this assumption implies that  $P$  is geodesically (strong) convex and smooth on  $\mathcal{B}$ . Under Assumption 5.2, the global minimum exists and is unique if all  $\{\phi_{k,t}\}_{k=1}^K$  are contained in  $\mathcal{B}$ , i.e.,  $P : \mathcal{B} \rightarrow \mathcal{A}$  is well-defined (Tron et al., 2012).

Recall the following trigonometric distance bound essential in the Riemannian optimization analysis.

**Lemma 5.3.** (Bonnabel, 2013; Zhang & Sra, 2016) If  $a, b, c$  are the side lengths of a geodesic triangle in a Riemannian manifold with sectional curvature lower bounded by  $\kappa_{\min}$ , and  $A$  is the angle between sides  $b$  and  $c$  (defined through the inverse exponential map and inner product in tangent space), then

$$a^2 \leq \frac{\sqrt{|\kappa_{\min}|}c}{\tanh(\sqrt{|\kappa_{\min}|}c)}b^2 + c^2 - 2bc \cos(A). \quad (9)$$

We define the following key geometric constant that captures the impact of manifold curvature:

$$\zeta = \begin{cases} \frac{\sqrt{|\kappa_{\min}|}D}{\tanh(\sqrt{|\kappa_{\min}|}D)}, & \text{if } \kappa_{\min} < 0, \\ 1, & \text{if } \kappa_{\min} \geq 0, \end{cases} \quad (10)$$

Note that most (if not all) practical manifold optimization problems can satisfy these assumptions.

Leveraging Assumption 5.1 and Lemma 5.3, we can readily establish the following corollary.

**Corollary 5.4.** (Zhang & Sra, 2016) For any Riemannian manifold  $\mathcal{M}$  where the sectional curvature is lower bounded by  $\kappa_{\min}$  and for any points  $x, x_t \in \mathcal{M}$ , the update  $x_{t+1} = \exp_{x_t}(-\mu \nabla F(x_t))$  satisfies the inequality:

$$\begin{aligned}\langle -\nabla F(x_t), \exp_{x_t}^{-1}(x) \rangle &\leq \frac{1}{2\mu} (d^2(x_t, x) - d^2(x_{t+1}, x)) \\ &\quad + \frac{\zeta\mu}{2} \|\nabla F(x_t)\|^2. \end{aligned} \quad (11)$$

This corollary unveils a significant relationship between two consecutive updates within an iterative optimization algorithm on a manifold with curvature bounded from below.

Part of our analysis will be performed under the following assumption of geodesically convex risk functions.

**Assumption 5.5 (Geodesical convexity).** A function  $J_k : \mathcal{M} \rightarrow \mathbb{R}$  is geodesically convex (g-convex) if for any  $x, y \in$



$\mathcal{M}$ , a geodesic  $\gamma$  such that  $\gamma(0) = x$  and  $\gamma(1) = y$ , and  $\alpha \in [0, 1]$ , we have:

$$J_k(\gamma(\alpha)) \leq (1 - \alpha)J_k(x) + \alpha J_k(y), \quad (12)$$

or equivalently, we have

$$J_k(y) \geq J_k(x) + \langle \nabla J_k(x), \exp_x^{-1}(y) \rangle. \quad (13)$$

Meanwhile, we require the risk function  $J_k$  at each agent to be geodesically smooth.

**Assumption 5.6 (Geodesic smoothness).** A differentiable function  $J_k$  is geodesically  $L$ -smooth ( $L$ -g-smooth) if its gradient is  $L$ -Lipschitz, i.e., for any  $x, y \in \mathcal{M}$ , it satisfies:

$$J_k(y) \leq J_k(x) + \langle \nabla J_k(x), \exp_x^{-1}(y) \rangle + \frac{L}{2} \|\exp_x^{-1}(y)\|^2, \quad (14)$$

where the gradient of a function  $J_k : \mathcal{M} \rightarrow \mathbb{R}$  is said to be  $L$ -Lipschitz if, for any  $x, y \in \mathcal{M}$  in the domain of  $J_k$ , it satisfies:

$$\|\nabla J_k(x) - \Gamma_y^x \nabla J_k(y)\| \leq L \|\exp_x^{-1}(y)\|, \quad (15)$$

where  $\Gamma_y^x$  denotes the parallel transport operator from  $y$  to  $x$ .

In addition, we make assumptions about the average and second moment of the gradient noise process.

**Assumption 5.7 (Gradient noise process).** Denote  $\mathcal{F}_t$  as the filtration generated by the random process  $\mathbf{w}_{k,s}$  for all  $k$  and for  $s \leq t$ , that is,

$$\mathcal{F}_t \triangleq \{\mathbf{w}_0, \mathbf{w}_1, \dots, \mathbf{w}_t\}, \quad (16)$$

where  $\mathbf{w}_s \triangleq \text{col}\{\mathbf{w}_{1,s}, \dots, \mathbf{w}_{K,s}\}$  contains the iterates across the network at time  $s$ . Define  $S_{t+1}(\mathbf{w}_t) \triangleq \widehat{\nabla J}(\mathbf{w}_t) - \nabla J(\mathbf{w}_t)$  as the gradient noise process at the time instant  $t$ . It is assumed that

$$\mathbb{E}\{S_{t+1}(\mathbf{w}_t) | \mathcal{F}_t\} = \mathbf{0}, \quad (17)$$

$$\mathbb{E}\{\|S_{t+1}(\mathbf{w}_t)\|^2 | \mathcal{F}_t\} \leq \sigma^2, \quad (18)$$

for some non-negative constant  $\sigma$ .

With these assumptions, we can build some preliminary lemmas that will be used in the proof of our main results.

## 5.2. Preliminary lemmas

We first establish a lemma that bounds the gradient of the penalty function  $P$  in terms of the penalty itself.

**Lemma 5.8.** *Under Assumption 5.2, for the gradient of the penalty, it holds that*

$$\|\nabla P(\phi_t)\|^2 \leq 2P(\phi_t). \quad (19)$$

*Proof.* Appendix A.1.  $\square$

Under Assumption 5.5, we can establish the following property for the risk function  $J$ .

**Lemma 5.9.** *Under Assumption 5.5, define  $\bar{\mathbf{w}} = \text{col}\{\mathbf{w}_m, \dots, \mathbf{w}_m\}$  with  $\mathbf{w}_m$  being the Fréchet mean (barycenter) of  $\mathbf{w}_1, \dots, \mathbf{w}_K$ , we have*

$$J(\bar{\mathbf{w}}) \leq J(\mathbf{w}). \quad (20)$$

*Proof.* Appendix A.2.  $\square$

This proof follows similarly to Proposition 10 in (Yokota, 2016) and Theorem 1.1 in (Paris, 2020).

The following lemma builds on assumptions 5.1 and 5.6, and establishes an upper bound on  $\|\nabla J(\mathbf{w}_t)\|$ .

**Lemma 5.10.** *Under assumptions 5.1 and 5.6, we have:*

$$\|\nabla J(\mathbf{w}_t)\| \leq G, \quad (21)$$

for a non-negative constant  $G < \infty$ .

*Proof.* Appendix A.3.  $\square$

This upper bound is similar to the one used in (Shah, 2017; Deng & Hu, 2023) under a diminishing step size.

## 5.3. Network agreement

To begin with, we first show that the Riemannian diffusion adaptation algorithm approximately converges toward network agreement. In other words,  $\mathbf{w}_t$  converges to the consensus submanifold  $\mathcal{A}$  with high probability. The following lemma builds on Lemma 5.10 under the additional conditions set forth in assumption 5.7 and 5.2.

**Lemma 5.11.** *Under assumptions 5.1, 5.2, 5.6 and 5.7, suppose  $\alpha \in (0, h_{\max}^{-1}]$ . The sequence  $\{P(\phi_t)\}_{t \geq 0}$  satisfies the following relation:*

$$\begin{aligned} \mathbb{E}\{P(\phi_{t+1}) - P(\phi_t)\} &\leq -\frac{\alpha}{4} \mathbb{E}\|\nabla P(\phi_t)\|^2 + \frac{5\mu^2}{\alpha} G^2 \\ &\quad + \frac{\mu^2}{\alpha} \sigma^2. \end{aligned} \quad (22)$$

*Proof.* Appendix A.4.  $\square$

This lemma reveals the evolution of the difference  $\mathbb{E}\{P(\phi_{t+1}) - P(\phi_t)\}$  in the optimization process. The first term on the right-hand side of (22) is strictly negative and suggests a decrease in the expectation of penalty by a magnitude proportional to  $\mathbb{E}\|\nabla P(\phi_t)\|^2$ . However, the second and third terms on the right-hand side of (22) could be large enough to allow the objective value to increase. In

the following, with an additional assumption that the cost  $J$  is geodesically convex (Assumption 5.5), we prove that the expectation of penalty decreases strictly and can be upper bounded with a small value after sufficient iterations.

**Theorem 5.12.** *Under assumptions 5.1, 5.2, 5.5, 5.6, and 5.7, suppose  $\alpha \in (0, h_{max}^{-1}]$ . The sequence  $\{P(\phi_t)\}_{t \geq 0}$  satisfies the following relation:*

$$\mathbb{E}\{P(\phi_t)\} \leq \frac{11\mu^2}{2\alpha\tau}G^2 + \frac{3\mu^2}{\alpha\tau}\sigma^2, \quad (23)$$

after sufficient iterations  $s_o$ , given by

$$s_o = \frac{2\log(\mu)}{\log(1-\tau)} + O(1) = O(\mu^{-1}), \quad (24)$$

where  $\tau = \min\{\frac{1}{2\zeta}, \alpha h_{min}\}$ ,  $O(1)$  denotes a constant term, and  $O(\mu^{-1})$  denotes a term that is equal or higher in order than  $\mu^{-1}$ , the last equality holds for sufficiently small  $\mu$ .

*Proof.* Appendix B.1.  $\square$

The result in Theorem 5.12 establishes that after sufficient iterations  $s_o = O(\mu^{-1})$ , we have:

$$\mathbb{E}\{P(\phi_t)\} \leq O(\mu^2), \quad (25)$$

or, from Markov's inequality:

$$\Pr\{P(\phi_t) \geq \mu\} \leq O(\mu), \quad (26)$$

which means the local estimates in  $\phi_t$  coalesce around  $\phi_t^* \in \mathcal{A}$  (where  $P(\phi_t^*) = 0$ ) with high probability. These results are consistent with Theorem 1 in (Vlaski & Sayed, 2021) where the Euclidean diffusion adaptation algorithm is analyzed for non-convex environments.

Combined with Lemma 5.8, Theorem 5.12 leads to the following corollary.

**Corollary 5.13.** *With the assumptions in Theorem 5.12. The sequence  $\{\|\nabla P(\phi_t)\|^2\}_{t \geq 0}$  satisfies the following relation:*

$$\mathbb{E}\|\nabla P(\phi_t)\|^2 \leq \frac{11\mu^2}{\alpha\tau}G^2 + \frac{6\mu^2}{\alpha\tau}\sigma^2, \quad (27)$$

after sufficient iterations  $s_o = O(\mu^{-1})$ .

Hence, according to  $\nabla P(\phi_t^*) = \mathbf{0}$  and the update in (7), we conclude that  $w_t$  approximately approaches  $\phi_t$  and achieves network agreement, or equivalently  $w_t \in \mathcal{A}$  with high probability after sufficient iterations.

#### 5.4. Non-asymptotic convergence

Next, we examine the convergence of Algorithm 1 after sufficient iterations  $s_o$ . For this purpose, we make use of the upper bound on  $\mathbb{E}\|\nabla P(\phi_t)\|^2$  given in Corollary 5.13. Before this, we introduce the following lemma that builds on the same assumptions as in Lemma 5.11 and can be regarded as a symmetric result of the relation (22).

**Lemma 5.14.** *Under assumptions 5.1, 5.2, 5.6 and 5.7, suppose  $\mu \in (0, L^{-1}]$ . The sequence  $\{J(w_t)\}_{t \geq 0}$  satisfies the following relation:*

$$\mathbb{E}\{J(w_{t+1}) - J(w_t)\} \leq -\frac{\mu}{4}\mathbb{E}\|\widehat{\nabla J}(w_t)\|^2 + \frac{5\alpha^2}{\mu}\mathbb{E}\|\nabla P(\phi_{t+1})\|^2. \quad (28)$$

*Proof.* Appendix A.5.  $\square$

This lemma shows the evolution of the term  $\mathbb{E}\{J(w_{t+1}) - J(w_t)\}$  in the optimization process. The strictly negative term on the right-hand side of (28) suggests a decrease in the expectation of the risk function by a magnitude proportional to  $\mathbb{E}\|\widehat{\nabla J}(w_t)\|^2$ , while the positive one could be large enough to allow the objective value to increase.

In the following, we consider (and study the convergence of) a streaming average of iterates  $\{w_{s_o+1}, \dots, w_t\}$ , given by  $\{w'_{s_o+1}, \dots, w'_t\}$  with  $w'_{s_o+1} = w_{s_o+1}$ ,  $w'_{s+1} = \exp_{w'_s} \left( \frac{1}{s-s_o+1} \exp_{w'_s}^{-1}(w_{s+1}) \right)$  for  $s_o + 1 \leq s \leq t-2$ , and

$$w'_t = \exp_{w'_{t-1}} \left( \frac{2\zeta}{2\zeta + t - s_o - 1} \exp_{w'_{t-1}}^{-1}(w_t) \right). \quad (29)$$

This provides a natural way of averaging along a trajectory restricted to a manifold (Tripuraneni et al., 2018). For example, when  $\mathcal{M}$  is a Euclidean space, we can write  $\exp_x(v)$  as  $x + v$ , and the streaming average reduces to  $w'_{s_o+1} = w_{s_o+1}$ ,  $w'_{s+1} = w'_s + \frac{1}{s-s_o+1}(w'_s - w_{s+1})$  for  $s_o + 1 \leq s \leq t-2$  and  $w'_t = w'_{t-1} + \frac{2\zeta}{2\zeta + t - s_o - 1}(w'_{t-1} - w_t)$ . Inspired by (Zhang & Sra, 2016), we design a Lyapunov function of  $w_t$  as

$$\Delta'_t \triangleq J(w'_t) - J(w^*), \quad (30)$$

with auxiliary variables  $w'_t$  defined in (29) and  $w^*$  denoted as the optimal solution to (2). Under an additional assumption of geodesic convexity (Assumption 5.5), we can establish the following result that  $\mathbb{E}\Delta'_t$  decreases strictly and can be bounded above.

**Theorem 5.15.** *Under assumptions 5.1, 5.2, 5.5, 5.6 and 5.7, suppose  $\alpha \in (0, h_{max}^{-1}]$  and  $\mu \in (0, L^{-1}]$ . The sequence  $\{J(w'_t)\}_{t \geq s_o+1}$  satisfies the following relation:*

$$\mathbb{E}\Delta'_t \leq \frac{\zeta LD^2 + (t - s_o) \left( \frac{231\zeta\alpha\mu}{2\tau}G^2 + \frac{63\zeta\alpha\mu}{\tau}\sigma^2 \right)}{2\zeta + t - s_o - 1}. \quad (31)$$

*Proof.* Appendix B.2.  $\square$

Theorem 5.15 establishes that non-asymptotic convergence of Algorithm 1 can be guaranteed after sufficient iterations

for sufficiently small step sizes  $\mu$  and  $\alpha$ , if  $J$  is geodesically convex and smooth.

Compared to the Euclidean counterpart (Chen & Sayed, 2012; Sayed et al., 2013; Vlaski & Sayed, 2021), key differences in our analysis include the impact of manifold curvature  $\kappa$  (captured in the parameter  $\zeta$ ) and the non-linear nature of the combination step (4). This makes traditional techniques like adjacency matrix decomposition unfeasible, since the network centroid cannot be computed using simple linear expressions. We address these challenges through a novel framework that studies network agreement via the evolution of the penalty term  $P(\phi_t)$ , and establish non-asymptotic convergence results using the carefully designed Lyapunov function in (30).

## 6. Examples and applications

In this section, we tailor Algorithm 1 to two common instances of Riemannian manifolds. The first one is the Grassmann manifold, a set of  $k$ -dimensional linear subspaces of  $\mathbb{R}^p$ , denoted by  $\mathcal{G}_n^p$ . The second is the manifold of  $p \times p$  SPD matrices, denoted by  $\mathcal{S}_n^{++}$ . We consider applying our algorithms on  $\mathcal{G}_n^p$  and a product manifold involving  $\mathcal{S}_n^{++}$  to online distributed PCA and GMM inference, respectively. While the exponential map is convenient for theoretical analysis, the retractions often lead to more practical and efficient computations. Thus, for computational simplicity, we replace the exponential maps in the updates (3) and (4) with approximate retractions as in (Bonnabel, 2013). We provide definitions of the geodesic distance, Riemannian gradient, and retraction of  $\mathcal{G}_n^p$  and  $\mathcal{S}_n^{++}$  in Appendix C. The computational complexity of our algorithm is discussed in Appendix E.

### 6.1. Distributed PCA

We consider applying our algorithm on  $\mathcal{G}_n^p$  to the online distributed PCA problem with  $\mathbf{x}_k \in \mathbb{R}^n$  being data samples observed by each agent  $k$ . In the decentralized setting, we consider the following problem:

$$\min_{\pi(\mathbf{U}_k) \in \mathcal{G}_n^p} -\mathbb{E}_{\mathbf{x}_k} \{ \text{tr}(\mathbf{U}_k^T \mathbf{x}_k \mathbf{x}_k^T \mathbf{U}_k) \}, \quad (32)$$

where  $\pi(\mathbf{U}_k)$  represents the local estimate at agent  $k$ . The expectation in the loss function (32) is approximated by realizations  $\mathbf{x}_{k,t}$  at each time instant  $t$ . Note that although various works formulate PCA on the Stiefel manifold (Chen et al., 2021; Wang & Liu, 2022; Wang et al., 2023), the loss function in (32) is invariant to orthonormal transformations. Thus, we formulate the problem on the Grassmannian manifold since it makes the solution unique (Cunningham & Ghahramani, 2015). This formulation has also been found to have a similar mathematical structure of strong geodesic convexity, allowing arguments from convex optimization on

manifolds to be applied (Alimisis & Vandereycken, 2024). The Riemannian stochastic gradient of the loss function in (32) on  $\mathcal{G}_n^p$  is computed using the Euclidean gradient of (32) at  $\mathbf{U}_{k,t}$  and (90) given in Appendix C.1, leading to:

$$h(\mathbf{U}_{k,t}, \mathbf{x}_{k,t}) = 2(\mathbf{I} - \mathbf{U}_{k,t} \mathbf{U}_{k,t}^T) \mathbf{x}_{k,t} \mathbf{x}_{k,t}^T \mathbf{U}_{k,t}.$$

The retraction used is defined in (91). In order to evaluate the accuracy of the solutions, we consider the geodesic distance (89) between the estimates at each time instant  $\pi(\mathbf{U}_{k,t})$  and the optimal solution  $\pi(\mathbf{U}^*)$ , and we define the mean square deviation (MSD) accordingly as

$$\text{MSD} = \frac{1}{K} \sum_{k=1}^K d_{\mathcal{G}_n^p}^2(\mathbf{U}_{k,t}, \mathbf{U}^*).$$

### 6.2. Distributed GMM inference

Another challenging application of our algorithm is distributed parameter estimation for GMMs with  $\mathbf{x}_k \in \mathbb{R}^n$  being data samples observed by each agent  $k$ . The decentralized inference of mixtures of  $M$  Gaussians with coefficients  $\boldsymbol{\rho} \triangleq \{\rho_1, \dots, \rho_M\}$ , whose probability density is  $p(\mathbf{x}) \triangleq \sum_{i=1}^M \rho_i p_{\mathcal{N}}(\mathbf{x}; \mathbf{m}_i, \boldsymbol{\Sigma}_i)$  with  $p_{\mathcal{N}}$  a multivariate Gaussian with mean  $\mathbf{m}_i$  and covariance  $\boldsymbol{\Sigma}_i \succ 0$ , can be reformulated as in (Hosseini & Sra, 2015):

$$\min_{\substack{\{\mathbf{S}_i\}_{i=1}^M \\ \{\eta_i\}_{i=1}^{M-1}}} -\mathbb{E}_{\mathbf{x}_k} \left\{ \log \left( \sum_{i=1}^M \frac{e^{\eta_i}}{\sum_{i=1}^M e^{\eta_i}} q_{\mathcal{N}}(\mathbf{y}_k; \mathbf{S}_i) \right) \right\}, \quad (33)$$

where  $\mathbf{y}_k^T = [\mathbf{x}_k^T \ 1]$ ,  $\eta_i = \log \frac{\rho_i}{\rho_M}$  for  $i = 1, \dots, M-1$  and  $\eta_M = 0$ , which makes the problem unconstrained (Jordan & Jacobs, 1994), and  $q_{\mathcal{N}}(\mathbf{y}_k; \mathbf{S}_i) = \sqrt{2\pi} e^{\frac{1}{2}} p_{\mathcal{N}}(\mathbf{y}_k; \mathbf{0}, \mathbf{S}_i)$ . The problem (33) reformulated on the product manifold  $\prod_{i=1}^M \mathcal{S}_n^{++} \times \mathbb{R}^{M-1}$  has the same optimum as that of the original log-likelihood of  $p(\mathbf{x})$  (Hosseini & Sra, 2015), i.e.,

$$\mathbf{S}_i^* = \begin{pmatrix} \boldsymbol{\Sigma}_i^* + \mathbf{m}_i^* \mathbf{m}_i^{*T} & \mathbf{m}_i^* \\ \mathbf{m}_i^{*T} & 1 \end{pmatrix}.$$

The log-likelihood has been shown to be geodesically convex for the case of a single Gaussian (Hosseini & Sra, 2015), but not necessarily for multiple Gaussians. From this example, we can find the proposed algorithm itself can work even in some situations when not all these assumptions are satisfied. The Riemannian gradient of the loss function in (33) on the product manifold is composed of the (Riemannian) gradients w.r.t.  $\{\mathbf{S}_i\}_{i=1}^M$  on  $\prod_{i=1}^M \mathcal{S}_n^{++}$  and  $\{\eta_i\}_{i=1}^{M-1}$  in  $\mathbb{R}^{M-1}$ . Specifically, the Riemannian gradient w.r.t.  $\mathbf{S}_i$  was computed via the Euclidean gradient of (33) at  $\mathbf{S}_{i,k,t}$  and (93) given in Appendix C.2. The retraction used is defined in (94). In this task, it is not very meaningful to compute the MSD values according to (92), because GMM

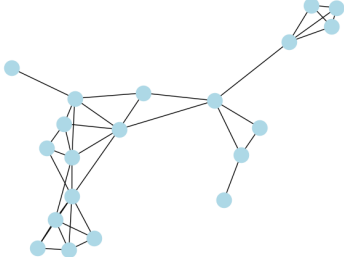


Figure 1. Graph topology.

is not inherently identifiable, which means that the parameters of the model may not be uniquely determined. Thus, to evaluate the performance of the solutions, we consider the average log-likelihood (ALL) as in (Hosseini & Sra, 2015; Collas et al., 2023).

## 7. Numerical experiments

In this section, we present numerical experiments on distributed PCA and parameter estimation of GMMs, which are formulated on manifolds as explained in Section 6. Our method is implemented in Python with the Pymanopt toolbox (Townsend et al., 2016). Open-source code to reproduce the results is publicly available on [https://github.com/xiuheng-wang/diffusion\\_manifold\\_release](https://github.com/xiuheng-wang/diffusion_manifold_release). The graph topology of the multi-agent system used for the experiments is illustrated in Figure 1. The weights in matrix  $C$  were randomly generated by the Metropolis rule (Xiao et al., 2006) with  $K = 20$  agents<sup>3</sup>. For simulation on synthetic data, the MSD results are averaged over 100 times independent Monte Carlo experiments. Hereafter, we briefly describe the baselines.

**Baselines:** We compare our algorithm against the Riemannian non-cooperative and centralized strategies for both PCA and GMM inference. The non-cooperative algorithm independently applies R-SGD on each agent using its local data  $\mathbf{x}_{k,t}$ , while the centralized works on data  $\mathbf{X}_t = \{\mathbf{x}_{k,t}\}_{k=1}^K$  collected from all agents. We also provide comparisons with an extrinsic consensus algorithm on the Stiefel manifold: Decentralized Riemannian Stochastic Gradient Descent (DRSGD) (Chen et al., 2021) for PCA. For GMM inference, to the best of our knowledge there are no approaches that are both online and decentralized. Thus, for comparison, we extend the decentralized consensus SGD (Nedic et al., 2010; Lian et al., 2017) to the product manifold presented in Section 6.2 using a projection operator to ensure the constraints are satisfied. This Extrinsic Consensus strategy for GMM inference is named ECGMM.

<sup>3</sup>To illustrate the applicability to more networks, we include additional experimental results in Appendix D.2.

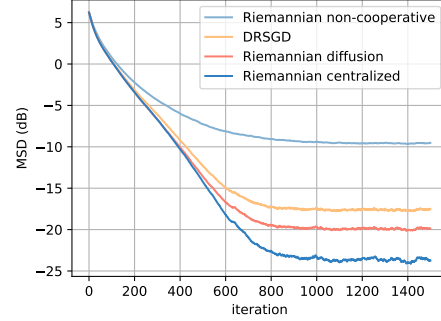


Figure 2. Illustration of MSD performance of the algorithms for distributed PCA on synthetic data.

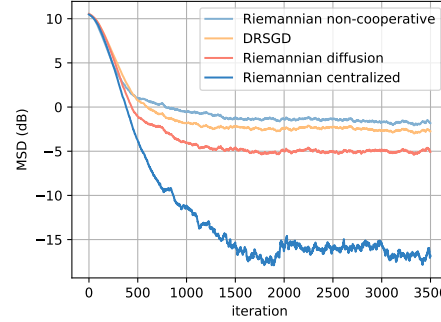


Figure 3. Illustration of MSD performance of the algorithms for distributed PCA on real data.

### 7.1. Experiments on PCA

We first present results on PCA formulated on  $\mathcal{G}_n^p$  with both synthetic and real data.

**Synthetic data:** We generate synthetic data as in (Chen et al., 2021). First, we set  $n = 10$ ,  $p = 5$ , and independently sample  $1500K$  data points according to a multivariate Gaussian model to obtain a matrix  $\mathbf{S} \in \mathbb{R}^{n \times 1500K}$ . Let  $\mathbf{S} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$  be its truncated SVD. We modify the distribution of  $\mathbf{\Lambda}$  as  $\mathbf{\Lambda}' = \text{diag}(\lambda^i)$  with  $\lambda = 0.8$  and  $i = 0, \dots, n-1$  to reset  $\mathbf{S}$  as  $\mathbf{S}' = \mathbf{U}\mathbf{\Lambda}'\mathbf{V}^T$ . We randomly shuffle and split the columns of  $\mathbf{S}' \in \mathbb{R}^{n \times 1500K}$  into 1500 subsets to obtain  $\mathbf{X}_t$  for all time instants  $t = 1, \dots, 1500$ . The simulations used fixed step sizes  $\mu = 0.05$  and  $\alpha = 0.8$ . For our algorithm, the step sizes control the tradeoff between convergence speed and steady-state performance; this is illustrated with experimental results in Appendix D.3.

**Real data:** We also obtain numerical results on the MNIST dataset (LeCun, 1998). The dataset contains 70000 hand-written images with  $n = 784$  pixels. The data matrix is normalized such that the elements are in the range  $[0, 1]$  and then centered. To compute MSD, we perform PCA on the full data matrix and regard its result as the optimum. We randomly shuffle the images, partition them into  $K = 20$  subsets, and then run the algorithms to compute the first  $p = 5$  principal components with the fixed step sizes  $\mu = 0.002$  and  $\alpha = 0.005$ .



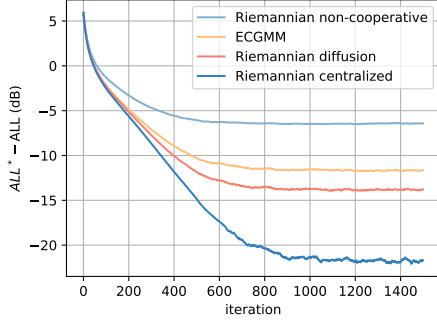


Figure 4. Illustration of ALL differences of the algorithms for distributed GMM inference on synthetic data.

**Discussion:** Figure 2 shows the MSD learning curves for the compared algorithms on synthetic data. It can be seen that the Riemannian diffusion adaptation strategy achieves a significant improvement in MSD performance compared to the non-cooperative case, which indicates the benefit of information exchange. Moreover, our method also outperforms DRSGD. The centralized solution achieves the lowest MSD, as it can access information over the whole graph. The proposed algorithm is fully decentralized, where each agent uses only locally observed data to update its local estimate and exchange information only among neighboring agents. Although the proposed algorithm has lower performance compared to the centralized method, it can be computed in parallel on multiple agents. The MSD of the different methods on real data, shown in Figure 3, behaves similarly to that in the experiment with synthetic data, showing the same comparative performances between the different approaches.

## 7.2. Experiments on GMM inference

Now we show results on a more challenging task: GMM inference formulated on  $\prod_{i=1}^M \mathcal{S}_n^{++} \times \mathbb{R}^{M-1}$  with synthetic and real data. As in (Hosseini & Sra, 2015), we initialize the mixture parameters for all the methods using k-means++.

**Synthetic data:** To generate synthetic data, we choose the parameters  $\mathbf{m}$ ,  $\Sigma$  and  $\rho$  of the Gaussian mixture similarly to (Collas et al., 2023). First,  $\Sigma$  is generated using its eigen-decomposition  $\Sigma = U\Lambda U^T$ .  $U$  is drawn from the uniform distribution on the orthogonal group, and all the elements of  $\Lambda$  are drawn from a chi-squared distribution with a degree of freedom 3. Second, the elements of  $\rho$  are drawn from a Gamma distribution with a shape parameter 10, and then normalized. Third,  $\mathbf{m}$  is sampled from a multivariate Gaussian distribution  $\mathcal{N}(\mathbf{0}, 3\mathbf{I})$ , where sampling is repeated until the following inequality is satisfied (Hosseini & Sra, 2015):

$$\forall i, j, \|\mathbf{m}_i - \mathbf{m}_j\| \geq \max\{tr(\Sigma_i), tr(\Sigma_j)\}.$$

We set  $n = 8$  and  $M = 3$  and independently sample 1500K data points according to the Gaussian mixture above. These

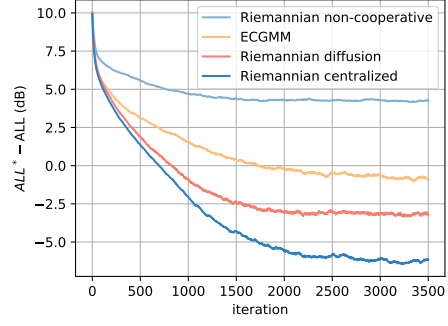


Figure 5. Illustration of ALL differences of the algorithms for distributed GMM inference on real data.

data points are randomly shuffled and split into 1500 subsets to obtain  $\mathbf{X}_t$  for all time instants  $t = 1, \dots, 1500$ . The simulations used fixed step sizes  $\mu = 0.04$  and  $\alpha = 0.05$ .

**Real data:** Again, we perform the same data processing for the MNIST dataset (LeCun, 1998) as in Section 7.1. Then, we apply PCA to reduce the dimensionality  $n = 20$ . We compute the Expectation Maximization (EM) solution on the full dataset and regard its result as an optimum. To evaluate the performance, we compare the difference between the ALL values of the optimum and estimated solutions, denoted as  $ALL^* - ALL$ . We implement the compared algorithms to infer mixtures of Gaussian models with  $M = 7$  and fixed step sizes  $\mu = 0.08$  and  $\alpha = 0.08$ .

**Discussion:** Figure 4 and Figure 5 illustrate the ALL difference values for the compared methods on synthetic and real data, respectively. The results demonstrate that our method outperforms both the non-cooperative algorithm and ECGMM, further highlighting its effectiveness. As expected, the centralized case achieves the lowest ALL difference.

## 8. Conclusions

In this paper, the Riemannian diffusion adaptation algorithm is proposed. The strategy consists of two efficient steps: an adaptation step, where R-SGD is used at each agent to update the estimate of the local solution on the manifold, and a combination step, where the estimates of neighboring agents are combined on the tangent space. A theoretical analysis is provided under constant step size, showing that network agreement is achieved with high probability and the algorithm converges non-asymptotically to a neighborhood of the optimal solution. The proposed method is applied to on-line decentralized PCA and GMM inference. Experimental results on both synthetic and real-world data illustrate the efficacy of the proposed strategy. One main limitation of this work is that the theoretical results rely on the use of the exponential map, which can be computationally heavy. This is discussed in more detail in Appendix F.

## Impact Statements

This paper presents work that aims to advance the field of Machine Learning. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

## Acknowledgements

The authors would like to thank the reviewers for their constructive feedback. The work of Cédric Richard was supported in part by the French Government through the 3IA Côte d’Azur Investments in the Future Project under grant ANR-19-P3IA-0002, and in part by grant ANR-19-CE48-0002. The work of Ricardo Borsoi was supported in part by the French National Research Agency, under grants ANR-23-CE23-0024, ANR-23-CE94-0001, and by the National Science Foundation, under grant NSF 2316420. Xiuheng Wang would like to thank Dr. Mengfei Zhang for the beneficial discussion in the early exploratory stage of this work.

## References

- Absil, P.-A., Mahony, R., and Sepulchre, R. *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, 2009.
- Afsari, B. Riemannian  $\ell^p$  center of mass: existence, uniqueness, and convexity. *Proceedings of the American Mathematical Society*, 139(2):655–673, 2011.
- Afsari, B., Tron, R., and Vidal, R. On the convergence of gradient descent for finding the Riemannian center of mass. *SIAM Journal on Control and Optimization*, 51(3): 2230–2260, 2013.
- Alimisis, F. and Vandereycken, B. Geodesic convexity of the symmetric eigenvalue problem and convergence of steepest descent. *Journal of Optimization Theory and Applications*, pp. 1–40, 2024.
- Bianchi, P. and Jakubowicz, J. Convergence of a multi-agent projected stochastic gradient algorithm for non-convex optimization. *IEEE Transactions on Automatic Control*, 58(2):391–405, 2012.
- Blatt, D., Hero, A. O., and Gauchman, H. A convergent incremental gradient method with a constant step size. *SIAM Journal on Optimization*, 18(1):29–51, 2007.
- Bonnabel, S. Stochastic gradient descent on Riemannian manifolds. *IEEE Transactions on Automatic Control*, 58(9):2217–2229, 2013.
- Boumal, N. *An Introduction to Optimization on Smooth Manifolds*. Cambridge University Press, 2023.
- Boumal, N. and Absil, P.-a. Rtrmc: A Riemannian trust-region method for low-rank matrix completion. *Advances in Neural Information Processing Systems*, 24, 2011.
- Chen, J. and Sayed, A. H. Diffusion adaptation strategies for distributed optimization and learning over networks. *IEEE Transactions on Signal Processing*, 60(8):4289–4305, 2012.
- Chen, J. and Sayed, A. H. On the learning behavior of adaptive networks—part I: Transient analysis. *IEEE Transactions on Information Theory*, 61(6):3487–3517, 2015.
- Chen, J., Ye, H., Wang, M., Huang, T., Dai, G., Tsang, I., and Liu, Y. Decentralized Riemannian conjugate gradient method on the stiefel manifold. In *The Twelfth International Conference on Learning Representations*, 2024.
- Chen, S., Garcia, A., Hong, M., and Shahrampour, S. Decentralized Riemannian gradient descent on the Stiefel manifold. In *International Conference on Machine Learning*, pp. 1594–1605. PMLR, 2021.
- Collas, A., Breloy, A., Ren, C., Ginolhac, G., and Ovarlez, J.-P. Riemannian optimization for non-centered mixture of scaled gaussian distributions. *IEEE Transactions on Signal Processing*, 2023.
- Cunningham, J. P. and Ghahramani, Z. Linear dimensionality reduction: Survey, insights, and generalizations. *The Journal of Machine Learning Research*, 16(1):2859–2900, 2015.
- Deng, K. and Hu, J. Decentralized projected Riemannian gradient method for smooth optimization on compact submanifolds. *arXiv:2304.08241*, 2023.
- Di Lorenzo, P. and Scutari, G. Next: In-network nonconvex optimization. *IEEE Transactions on Signal and Information Processing over Networks*, 2(2):120–136, 2016.
- Edelman, A., Arias, T. A., and Smith, S. T. The geometry of algorithms with orthogonality constraints. *SIAM journal on Matrix Analysis and Applications*, 20(2):303–353, 1998.
- Hosseini, R. and Sra, S. Matrix manifold optimization for gaussian mixtures. *Advances in Neural Information Processing Systems*, 28, 2015.
- Hsieh, Y.-P., Karimi Jaghargh, M. R., Krause, A., and Mertikopoulos, P. Riemannian stochastic optimization methods avoid strict saddle points. *Advances in Neural Information Processing Systems*, 36, 2024.
- Huang, L.-K. and Pan, S. Communication-efficient distributed PCA by Riemannian optimization. In *International Conference on Machine Learning*, pp. 4465–4474. PMLR, 2020.

- Huang, Z., Huang, W., Jawanpuria, P., and Mishra, B. Federated learning on Riemannian manifolds with differential privacy. *arXiv preprint arXiv:2404.10029*, 2024a.
- Huang, Z., Huang, W., Jawanpuria, P., and Mishra, B. Riemannian federated learning via averaging gradient stream. *arXiv preprint arXiv:2409.07223*, 2024b.
- Jordan, M. I. and Jacobs, R. A. Hierarchical mixtures of experts and the EM algorithm. *Neural Computation*, 6(2):181–214, 1994.
- Kraisler, S., Talebi, S., and Mesbahi, M. Consensus on lie groups for the Riemannian center of mass. In *2023 62nd IEEE Conference on Decision and Control (CDC)*, pp. 4461–4466. IEEE, 2023a.
- Kraisler, S., Talebi, S., and Mesbahi, M. Distributed consensus on manifolds using the Riemannian center of mass. In *2023 IEEE Conference on Control Technology and Applications (CCTA)*, pp. 130–135. IEEE, 2023b.
- LeCun, Y. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- Lee, J. M. *Introduction to Smooth Manifolds*. Springer, 2nd edition, 2013.
- Li, J. and Ma, S. Federated learning on Riemannian manifolds. *Applied Set-Valued Analysis and Optimization*, 5(2), 2023.
- Lian, X., Zhang, C., Zhang, H., Hsieh, C.-J., Zhang, W., and Liu, J. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. *Advances in Neural Information Processing Systems*, 30, 2017.
- Nedic, A., Ozdaglar, A., and Parrilo, P. A. Constrained consensus and optimization in multi-agent networks. *IEEE Transactions on Automatic Control*, 55(4):922–938, 2010.
- Paris, Q. Jensen’s inequality in geodesic spaces with lower bounded curvature. *arXiv:2011.08597*, 2020.
- Pennec, X., Fillard, P., and Ayache, N. A Riemannian framework for tensor computing. *International Journal of Computer Vision*, 66(1):41–66, 2006.
- Sarlette, A. and Sepulchre, R. Consensus optimization on manifolds. *SIAM journal on Control and Optimization*, 48(1):56–76, 2009.
- Sato, H., Kasai, H., and Mishra, B. Riemannian stochastic variance reduced gradient algorithm with retraction and vector transport. *SIAM Journal on Optimization*, 29(2): 1444–1472, 2019.
- Sayed, A. H. Adaptive networks. *Proceedings of the IEEE*, 102(4):460–497, 2014.
- Sayed, A. H., Tu, S.-Y., Chen, J., Zhao, X., and Towfic, Z. J. Diffusion strategies for adaptation and learning over networks: an examination of distributed strategies and network behavior. *IEEE Signal Processing Magazine*, 30(3):155–171, 2013.
- Sayed, A. H. et al. Adaptation, learning, and optimization over networks. *Foundations and Trends® in Machine Learning*, 7(4-5):311–801, 2014.
- Shah, S. M. Distributed optimization on Riemannian manifolds for multi-agent networks. *arXiv:1711.11196*, 2017.
- Tatarenko, T. and Touri, B. Non-convex distributed optimization. *IEEE Transactions on Automatic Control*, 62(8):3744–3757, 2017.
- Townsend, J., Koep, N., and Weichwald, S. Pymanopt: A python toolbox for optimization on manifolds using automatic differentiation. *Journal of Machine Learning Research*, 17(137):1–5, 2016.
- Tripuraneni, N., Flammarion, N., Bach, F., and Jordan, M. I. Averaging stochastic gradient descent on Riemannian manifolds. In *Conference on Learning Theory*, pp. 650–687. PMLR, 2018.
- Tron, R., Afsari, B., and Vidal, R. Riemannian consensus for manifolds with bounded curvature. *IEEE Transactions on Automatic Control*, 58(4):921–934, 2012.
- Vandereycken, B. Low-rank matrix completion by Riemannian optimization. *SIAM Journal on Optimization*, 23(2): 1214–1236, 2013.
- Vlaski, S. and Sayed, A. H. Distributed learning in non-convex environments—Part I: Agreement at a linear rate. *IEEE Transactions on Signal Processing*, 69:1242–1256, 2021.
- Vlaski, S., Kar, S., Sayed, A. H., and Moura, J. M. Networked signal and information processing: Learning by multiagent systems. *IEEE Signal Processing Magazine*, 40(5):92–105, 2023.
- Vorontsov, E., Trabelsi, C., Kadoury, S., and Pal, C. On orthogonality and learning recurrent networks with long term dependencies. In *International Conference on Machine Learning*, pp. 3570–3578. PMLR, 2017.
- Wang, L. and Liu, X. Decentralized optimization over the Stiefel manifold by an approximate augmented Lagrangian function. *IEEE Transactions on Signal Processing*, 70:3029–3041, 2022.

- Wang, X., Jiao, Y., Wai, H.-T., and Gu, Y. Incremental aggregated Riemannian gradient method for distributed PCA. In *International Conference on Artificial Intelligence and Statistics*, pp. 7492–7510. PMLR, 2023.
- Wang, X., Borsoi, R. A., and Richard, C. Non-parametric online change point detection on Riemannian manifolds. In *International Conference on Machine Learning*, pp. 50143–50162. PMLR, 2024a.
- Wang, X., Borsoi, R. A., and Richard, C. Riemannian diffusion adaptation over graphs with application to online distributed PCA. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 9736–9740, 2024b.
- Wang, Y., Yin, W., and Zeng, J. Global convergence of ADMM in nonconvex nonsmooth optimization. *Journal of Scientific Computing*, 78:29–63, 2019.
- Weber, M. and Sra, S. Projection-free nonconvex stochastic optimization on Riemannian manifolds. *IMA Journal of Numerical Analysis*, 42(4):3241–3271, 2022.
- Xiao, L., Boyd, S., and Lall, S. A space-time diffusion scheme for peer-to-peer least-squares estimation. In *Proceedings of the 5th International Conference on Information Processing in Sensor Networks*, pp. 168–176, 2006.
- Yokota, T. Convex functions and barycenter on cat (1)-spaces of small radii. *Journal of the Mathematical Society of Japan*, 68(3):1297–1323, 2016.
- Yuan, K., Ying, B., Zhao, X., and Sayed, A. H. Exact diffusion for distributed optimization and learning—part I: Algorithm development. *IEEE Transactions on Signal Processing*, 67(3):708–723, 2018.
- Zeng, J. and Yin, W. On nonconvex decentralized gradient descent. *IEEE Transactions on Signal Processing*, 66(11):2834–2848, 2018.
- Zhang, H. and Sra, S. First-order methods for geodesically convex optimization. In *Conference on Learning Theory*, pp. 1617–1638, 2016.
- Zhang, H., Reddi, S. J., and Sra, S. Riemannian SVRG: Fast stochastic optimization on Riemannian manifolds. In *Advances in Neural Information Processing Systems*, pp. 4592–4600, 2016.



## A. Proofs of the lemmas

### A.1. Lemma 5.8

From the definition of  $P(\phi_t)$ , we have

$$\|\nabla P(\phi_t)\|^2 = \sum_{k=1}^K \left\| -\sum_{\ell=1}^K c_{\ell k} \exp_{\phi_{k,t}}^{-1}(\phi_{\ell,t}) \right\|^2 \leq \sum_{k=1}^K \sum_{\ell=1}^K c_{\ell k} \left\| \exp_{\phi_{k,t}}^{-1}(\phi_{\ell,t}) \right\|^2 = 2P(\phi_t), \quad (34)$$

where the inequality follows from Jensen's inequality and the fact that  $C$  is left-stochastic, as stated in Assumption 5.2.

### A.2. Lemma 5.9

Due to the  $g$ -convexity of  $J_k$  under Assumption 5.5, we have

$$J_k(\mathbf{w}_m) - J_k(\mathbf{w}_k) \leq \langle -\nabla J_k(\mathbf{w}_m), \exp_{\mathbf{w}_m}^{-1}(\mathbf{w}_k) \rangle \quad (35)$$

Multiply the above by  $\frac{1}{K}$  and sum over  $k$ , we get

$$\frac{1}{K} \sum_k J_k(\mathbf{w}_m) - \frac{1}{K} \sum_k J_k(\mathbf{w}_k) \leq \frac{1}{K} \sum_k \langle -\nabla J_k(\mathbf{w}_m), \exp_{\mathbf{w}_m}^{-1}(\mathbf{w}_k) \rangle \quad (36)$$

Using  $\frac{1}{K} \sum_k J_k(\mathbf{w}_m) := J(\bar{\mathbf{w}})$  and rearranging the terms, this can be rewritten as

$$J(\bar{\mathbf{w}}) - \frac{1}{K} \sum_k J_k(\mathbf{w}_k) \leq \langle -\nabla J_k(\mathbf{w}_m), \frac{1}{K} \sum_k \exp_{\mathbf{w}_m}^{-1}(\mathbf{w}_k) \rangle \quad (37)$$

Note that  $\frac{1}{K} \sum_k \exp_{\mathbf{w}_m}^{-1}(\mathbf{w}_k) = \frac{1}{K} \nabla_x \sum_k d^2(\mathbf{w}_k, \mathbf{x})|_{\mathbf{x}:=\mathbf{w}_m}$  is the gradient of the cost function  $\frac{1}{K} \sum_k d^2(\mathbf{w}_k, \mathbf{x})$  evaluated at the Fréchet mean  $\mathbf{w}_m$ , which is its minimizer, therefore, the first order optimality condition implies  $\frac{1}{K} \sum_k \exp_{\mathbf{w}_m}^{-1}(\mathbf{w}_k) = 0$ . Combining this with the previous results leads to  $J(\bar{\mathbf{w}}) \leq \frac{1}{K} \sum_k J_k(\mathbf{w}_k) = J(\mathbf{w})$ .

### A.3. Lemma 5.10

Since  $\mathcal{B}$  is compact (Assumption 5.1), from the geodesic smoothness of  $J_k$  (Assumption 5.6), we have:

$$\|\nabla J_k(\mathbf{w}_t)\| \leq G, \quad (38)$$

for a non-negative constant  $G < \infty$ . Observe that (38) implies a similar condition on the deviation from the centralized gradient via Jensen's inequality:

$$\|\nabla J(\mathbf{w}_t)\| = \left\| \frac{1}{K} \sum_k \nabla J_k(\mathbf{w}_t) \right\| \leq \frac{1}{K} \sum_k \|\nabla J_k(\mathbf{w}_t)\| \leq G. \quad (39)$$

### A.4. Lemma 5.11

Let us start with the update  $\mathbf{w}_t = \exp_{\phi_t}(-\alpha \nabla P(\phi_t))$  from (7), define  $\gamma_1(\alpha) \triangleq \exp_{\phi_t}(-\alpha \nabla P(\phi_t))$  as the minimal geodesic from  $\phi_t$  to  $\mathbf{w}_t$ , and use the second-order Taylor expansion of  $\alpha \mapsto P(\gamma_1(\alpha))$  around  $\alpha = 0$ , under assumptions 5.1 and 5.2, then we have (Tron et al., 2012)

$$\begin{aligned} P(\mathbf{w}_t) &\leq P(\phi_t) + \langle \nabla P(\phi_t), -\alpha \nabla P(\phi_t) \rangle + \frac{h_{max} \|\alpha \nabla P(\phi_t)\|^2}{2} \\ &= P(\phi_t) - \epsilon \|\nabla P(\phi_t)\|^2, \end{aligned} \quad (40)$$

where  $\epsilon \triangleq \alpha \left(1 - \frac{\alpha h_{max}}{2}\right) > 0$  since  $\alpha \in (0, h_{max}^{-1}]$ . Also, we use the first-order Taylor expansion of  $\alpha \mapsto \nabla P(\gamma_1(\alpha))$  around  $\alpha = 0$  to obtain the following bound:

$$\|\nabla P(\mathbf{w}_t) - \Gamma_{\phi_t}^{\mathbf{w}_t} \nabla P(\phi_t)\| \leq h_{max} \alpha \|\nabla P(\phi_t)\|. \quad (41)$$

Similarly, for the update  $\phi_{t+1} = \exp_{\mathbf{w}_t}(-\mu \widehat{\nabla J}(\mathbf{w}_t))$  from (6), define  $\gamma_2(\mu) \triangleq \exp_{\mathbf{w}_t}(-\mu \widehat{\nabla J}(\mathbf{w}_t))$  as the minimal geodesic from  $\mathbf{w}_t$  to  $\phi_{t+1}$ , use the second-order Taylor expansion of  $\mu \mapsto P(\gamma_2(\mu))$  around  $\mu = 0$ , under assumptions 5.1 and 5.2, then we have

$$P(\phi_{t+1}) \leq P(\mathbf{w}_t) + \langle \nabla P(\mathbf{w}_t), -\mu \widehat{\nabla J}(\mathbf{w}_t) \rangle + \frac{h_{max} \mathbb{E} \|\mu \widehat{\nabla J}(\mathbf{w}_t)\|^2}{2}. \quad (42)$$

Take the expectation on (42) w.r.t.  $\{\mathbf{x}_s\}_{s=0}^t$  and consider (17) and (18) in Assumption 5.7, we have

$$\begin{aligned} \mathbb{E}P(\phi_{t+1}) &\leq \mathbb{E}P(\mathbf{w}_t) + \mathbb{E}\{\langle \nabla P(\mathbf{w}_t), -\mu \widehat{\nabla J}(\mathbf{w}_t) \rangle\} + \frac{h_{max} \mathbb{E} \|\mu \widehat{\nabla J}(\mathbf{w}_t)\|^2}{2} \\ &= \mathbb{E}P(\mathbf{w}_t) + \mathbb{E}\{\langle \nabla P(\mathbf{w}_t), -\mu \mathbb{E}\{\widehat{\nabla J}(\mathbf{w}_t) | \mathcal{F}_t\} \rangle\} + \frac{h_{max} \mu^2}{2} \mathbb{E} \|\widehat{\nabla J}(\mathbf{w}_t)\|^2 \\ &= \mathbb{E}P(\mathbf{w}_t) + \mathbb{E}\{\langle \nabla P(\mathbf{w}_t), -\mu \nabla J(\mathbf{w}_t) \rangle\} + \frac{h_{max} \mu^2}{2} \mathbb{E} \|\widehat{\nabla J}(\mathbf{w}_t)\|^2 \\ &\leq \mathbb{E}P(\mathbf{w}_t) + \frac{\xi}{2} \mathbb{E} \|\nabla P(\mathbf{w}_t)\|^2 + \frac{1}{2\xi} \mu^2 \mathbb{E} \|\nabla J(\mathbf{w}_t)\|^2 + \frac{h_{max} \mu^2}{2} \mathbb{E} \|\widehat{\nabla J}(\mathbf{w}_t) - \nabla J(\mathbf{w}_t) + \nabla J(\mathbf{w}_t)\|^2 \\ &\leq \mathbb{E}P(\mathbf{w}_t) + \frac{\xi}{2} \mathbb{E} \|\nabla P(\mathbf{w}_t)\|^2 + \left(\frac{1}{2\xi} + h_{max}\right) \mu^2 \mathbb{E} \|\nabla J(\mathbf{w}_t)\|^2 + h_{max} \mu^2 \mathbb{E}\{\mathbb{E}\{\|\widehat{\nabla J}(\mathbf{w}_t) - \nabla J(\mathbf{w}_t)\|^2 | \mathcal{F}_t\}\} \\ &= \mathbb{E}P(\mathbf{w}_t) + \frac{\xi}{2} \mathbb{E} \|\nabla P(\mathbf{w}_t)\|^2 + \left(\frac{1}{2\xi} + h_{max}\right) \mu^2 \mathbb{E} \|\nabla J(\mathbf{w}_t)\|^2 + h_{max} \mu^2 \sigma^2, \end{aligned} \quad (43)$$

where we use the facts  $\langle a, b \rangle \leq \frac{\xi}{2} a^2 + \frac{1}{2\xi} b^2$  for  $\xi > 0$  and  $\frac{1}{2}(a+b)^2 \leq a^2 + b^2$  in the second and third inequalities, respectively. Next, we take the expectation on (40) w.r.t.  $\{\mathbf{x}_s\}_{s=0}^t$ , and combine the result with (43) to obtain

$$\mathbb{E}P(\phi_{t+1}) \leq \mathbb{E}P(\phi_t) - \epsilon \mathbb{E} \|\nabla P(\phi_t)\|^2 + \frac{\xi}{2} \mathbb{E} \|\nabla P(\mathbf{w}_t)\|^2 + \left(\frac{1}{2\xi} + h_{max}\right) \mu^2 \mathbb{E} \|\nabla J(\mathbf{w}_t)\|^2 + h_{max} \mu^2 \sigma^2. \quad (44)$$

Now we need to upper bound  $\mathbb{E} \|\nabla P(\mathbf{w}_t)\|^2$ . Consider

$$\begin{aligned} \frac{1}{2} \mathbb{E} \|\nabla P(\mathbf{w}_t)\|^2 &= \frac{1}{2} \mathbb{E} \|\nabla P(\mathbf{w}_t) - \Gamma_{\phi_t}^{\mathbf{w}_t} \nabla P(\phi_t) + \Gamma_{\phi_t}^{\mathbf{w}_t} \nabla P(\phi_t)\|^2 \\ &\leq \mathbb{E} \|\nabla P(\mathbf{w}_t) - \Gamma_{\phi_t}^{\mathbf{w}_t} \nabla P(\phi_t)\|^2 + \mathbb{E} \|\nabla P(\phi_t)\|^2 \\ &\leq (\alpha^2 h_{max}^2 + 1) \mathbb{E} \|\nabla P(\phi_t)\|^2, \end{aligned} \quad (45)$$

where the first inequality uses the fact  $\frac{1}{2}(a+b)^2 \leq a^2 + b^2$ , the second inequality uses (41). Plugging the upper bound of  $\frac{1}{2} \mathbb{E} \|\nabla P(\mathbf{w}_t)\|^2$ , as provided in (45), into (44) and reordering, we have

$$\begin{aligned} \mathbb{E}P(\phi_{t+1}) &\leq \mathbb{E}P(\phi_t) - \epsilon \mathbb{E} \|\nabla P(\phi_t)\|^2 + \xi(\alpha^2 h_{max}^2 + 1) \mathbb{E} \|\nabla P(\phi_t)\|^2 + \left(\frac{1}{2\xi} + h_{max}\right) \mu^2 \mathbb{E} \|\nabla J(\mathbf{w}_t)\|^2 + h_{max} \mu^2 \sigma^2 \\ &= \mathbb{E}P(\phi_t) - \frac{\epsilon}{2} \mathbb{E} \|\nabla P(\phi_t)\|^2 + \left(\frac{\alpha^2 h_{max}^2 + 1}{\epsilon} + h_{max}\right) \mu^2 \mathbb{E} \|\nabla J(\mathbf{w}_t)\|^2 + h_{max} \mu^2 \sigma^2 \\ &\leq \mathbb{E}P(\phi_t) - \frac{\epsilon}{2} \mathbb{E} \|\nabla P(\phi_t)\|^2 + \left(\frac{\alpha^2 h_{max}^2 + 1}{\epsilon} + h_{max}\right) \mu^2 G^2 + h_{max} \mu^2 \sigma^2, \end{aligned} \quad (46)$$

where in the equality we select  $\xi = \frac{\epsilon}{2(\alpha^2 h_{max}^2 + 1)}$  for simplicity, and in the second inequality we use (21) from Lemma 5.10. Since  $\alpha \in (0, h_{max}^{-1}]$ , we have  $h_{max} \leq \alpha^{-1}$  and  $\epsilon \geq \frac{\alpha}{2}$ , and thus we can further simplify (46) as

$$\mathbb{E}P(\phi_{t+1}) \leq \mathbb{E}P(\phi_t) - \frac{\alpha}{4} \mathbb{E} \|\nabla P(\phi_t)\|^2 + \frac{5\mu^2}{\alpha} G^2 + \frac{\mu^2}{\alpha} \sigma^2. \quad (47)$$

Re-arranging the terms in (47) gives the desired result.

### A.5. Lemma 5.14

Consider the smoothness property of  $J$  in Assumption 5.6 with  $\exp_{\mathbf{w}_t}^{-1}(\phi_{t+1}) = -\mu \widehat{\nabla J}(\mathbf{w}_t)$  from (6), we can write:

$$\begin{aligned} J(\phi_{t+1}) &\leq J(\mathbf{w}_t) + \langle \nabla J(\mathbf{w}_t), \exp_{\mathbf{w}_t}^{-1}(\phi_{t+1}) \rangle + \frac{L \|\exp_{\mathbf{w}_t}^{-1}(\phi_{t+1})\|^2}{2} \\ &= J(\mathbf{w}_t) + \langle \nabla J(\mathbf{w}_t), -\mu \widehat{\nabla J}(\mathbf{w}_t) \rangle + \frac{L \|\mu \widehat{\nabla J}(\mathbf{w}_t)\|^2}{2}. \end{aligned} \quad (48)$$

Also, we can obtain the following bound:

$$\|\nabla J(\phi_{t+1}) - \Gamma_{\mathbf{w}_t}^{\phi_{t+1}} \nabla J(\mathbf{w}_t)\| \leq L\mu \|\widehat{\nabla J}(\mathbf{w}_t)\|. \quad (49)$$

Take expectation on (48) w.r.t.  $\{\mathbf{x}_s\}_{s=0}^t$  and consider (17) in Assumption 5.7, we have:

$$\begin{aligned} \mathbb{E}J(\phi_{t+1}) &\leq \mathbb{E}J(\mathbf{w}_t) + \mathbb{E}\{\langle \nabla J(\mathbf{w}_t), -\mu \widehat{\nabla J}(\mathbf{w}_t) \rangle\} + \frac{L\mathbb{E}\|\mu \widehat{\nabla J}(\mathbf{w}_t)\|^2}{2} \\ &= \mathbb{E}J(\mathbf{w}_t) + \mathbb{E}\{\langle \mathbb{E}\{\widehat{\nabla J}(\mathbf{w}_t) | \mathcal{F}_t\}, -\mu \widehat{\nabla J}(\mathbf{w}_t) \rangle\} + \frac{L\mu^2}{2} \mathbb{E}\|\widehat{\nabla J}(\mathbf{w}_t)\|^2 \\ &= \mathbb{E}J(\mathbf{w}_t) - \epsilon \mathbb{E}\|\widehat{\nabla J}(\mathbf{w}_t)\|^2. \end{aligned} \quad (50)$$

where  $\epsilon \triangleq \mu \left(1 - \frac{\mu L}{2}\right) > 0$  since  $\mu \in (0, L^{-1}]$ . Again, consider the smoothness property of  $J$  in Assumption 5.6 with  $\exp_{\phi_{t+1}}^{-1}(\mathbf{w}_{t+1}) = -\alpha \nabla P(\phi_{t+1})$  from (7), we obtain:

$$\begin{aligned} J(\mathbf{w}_{t+1}) &\leq J(\phi_{t+1}) + \langle \nabla J(\phi_{t+1}), \exp_{\phi_{t+1}}^{-1}(\mathbf{w}_{t+1}) \rangle + \frac{L \|\exp_{\phi_{t+1}}^{-1}(\mathbf{w}_{t+1})\|^2}{2} \\ &= J(\phi_{t+1}) + \langle \nabla J(\phi_{t+1}), -\alpha \nabla P(\phi_{t+1}) \rangle + \frac{L \|\alpha \nabla P(\phi_{t+1})\|^2}{2} \\ &\leq J(\phi_{t+1}) + \frac{\xi}{2} \|\nabla J(\phi_{t+1})\|^2 + \left(\frac{1}{2\xi} + L\right) \alpha^2 \|\nabla P(\phi_{t+1})\|^2, \end{aligned} \quad (51)$$

where the second inequality uses the fact  $\langle a, b \rangle \leq \frac{\xi}{2} a^2 + \frac{1}{2\xi} b^2$ . Next, we take the expectation on (51) w.r.t.  $\{\mathbf{x}_s\}_{s=0}^t$ , and combine the result with (50) to obtain

$$\mathbb{E}J(\mathbf{w}_{t+1}) \leq \mathbb{E}J(\mathbf{w}_t) - \epsilon \mathbb{E}\|\widehat{\nabla J}(\mathbf{w}_t)\|^2 + \frac{\xi}{2} \mathbb{E}\|\nabla J(\phi_{t+1})\|^2 + \left(\frac{1}{2\xi} + L\right) \alpha^2 \mathbb{E}\|\nabla P(\phi_{t+1})\|^2, \quad (52)$$

Now we need to upper bound  $\mathbb{E}\|\nabla J(\phi_{t+1})\|^2$ . Consider

$$\begin{aligned} \frac{1}{2} \mathbb{E}\|\nabla J(\phi_{t+1})\|^2 &= \frac{1}{2} \mathbb{E}\|\nabla J(\phi_{t+1}) - \Gamma_{\mathbf{w}_t}^{\phi_{t+1}} \nabla J(\mathbf{w}_t) + \Gamma_{\mathbf{w}_t}^{\phi_{t+1}} \nabla J(\mathbf{w}_t)\|^2 \\ &\leq \mathbb{E}\|\nabla J(\phi_{t+1}) - \Gamma_{\mathbf{w}_t}^{\phi_{t+1}} \nabla J(\mathbf{w}_t)\|^2 + \mathbb{E}\|\nabla J(\mathbf{w}_t)\|^2 \\ &\leq (\mu^2 L^2 + 1) \mathbb{E}\|\widehat{\nabla J}(\mathbf{w}_t)\|^2, \end{aligned} \quad (53)$$

where the first inequality uses the fact  $\frac{1}{2}(a+b)^2 \leq a^2 + b^2$ , the second inequality uses (49) and the fact  $\mathbb{E}\|\nabla J(\mathbf{w}_t)\|^2 \leq \mathbb{E}\|\widehat{\nabla J}(\mathbf{w}_t)\|^2$ . Plugging the upper bound of  $\frac{1}{2} \mathbb{E}\|\nabla J(\phi_{t+1})\|^2$ , as provided in (53), into (52) and reordering, we have

$$\begin{aligned} \mathbb{E}J(\mathbf{w}_{t+1}) &\leq \mathbb{E}J(\mathbf{w}_t) - (\epsilon - \xi(\mu^2 L^2 + 1)) \mathbb{E}\|\widehat{\nabla J}(\mathbf{w}_t)\|^2 + \left(\frac{1}{2\xi} + L\right) \alpha^2 \mathbb{E}\|\nabla P(\phi_{t+1})\|^2 \\ &= \mathbb{E}J(\mathbf{w}_t) - \frac{\epsilon}{2} \mathbb{E}\|\widehat{\nabla J}(\mathbf{w}_t)\|^2 + \left(\frac{\mu^2 L^2 + 1}{\epsilon} + L\right) \alpha^2 \mathbb{E}\|\nabla P(\phi_{t+1})\|^2, \end{aligned} \quad (54)$$

where in the equality we select  $\xi = \frac{\epsilon}{2(\mu^2 L^2 + 1)}$  for simplicity. Since  $\mu \in (0, L^{-1}]$ , we have  $L \leq \mu^{-1}$  and  $\epsilon \geq \frac{\mu}{2}$ , and thus we can further simplify (54) as

$$\mathbb{E}J(\mathbf{w}_{t+1}) \leq \mathbb{E}J(\mathbf{w}_t) - \frac{\mu}{4} \mathbb{E}\|\widehat{\nabla J}(\mathbf{w}_t)\|^2 + \frac{5\alpha^2}{\mu} \mathbb{E}\|\nabla P(\phi_{t+1})\|^2. \quad (55)$$

Re-arranging the terms in (55) gives the desired result.

## B. Proofs of the theorems

### B.1. Theorem 5.12

Define  $\bar{\phi}_t = \text{col}\{\phi_{m,t}, \dots, \phi_{m,t}\}$  with  $\phi_{m,t}$  being the Fréchet mean of  $\phi_t$ , i.e.,  $\phi_{m,t} \triangleq \arg \min_{\phi} \sum_{k=1}^K d^2(\phi_{k,t}, \phi)$ . Further, define  $\gamma_3(\beta) \triangleq \exp_{\phi_t}(\beta \exp_{\phi_t}^{-1} \bar{\phi}_t)$  as the minimal geodesic from  $\phi_t$  to  $\bar{\phi}_t$ , use the second-order Taylor expansion of  $\beta \mapsto P(\gamma_3(\beta))$  around  $\beta = 0$ , then we have (Afsari et al., 2013):

$$\langle \nabla P(\phi_t), \exp_{\phi_t}^{-1}(\bar{\phi}_t) \rangle + \frac{h_{\min} \|\exp_{\phi_t}^{-1}(\bar{\phi}_t)\|^2}{2} \leq P(\bar{\phi}_t) - P(\phi_t), \quad (56)$$

Considering  $P(\bar{\phi}_t) = 0$ , we can further write

$$\begin{aligned} P(\phi_t) &= P(\phi_t) - P(\bar{\phi}_t) \leq \langle -\nabla P(\phi_t), \exp_{\phi_t}^{-1}(\bar{\phi}_t) \rangle - \frac{h_{\min} \|\exp_{\phi_t}^{-1}(\bar{\phi}_t)\|^2}{2} \\ &\leq \frac{1 - \alpha h_{\min}}{2\alpha} d^2(\phi_t, \bar{\phi}_t) - \frac{1}{2\alpha} d^2(\mathbf{w}_t, \bar{\phi}_t) + \frac{\zeta\alpha}{2} \|\nabla P(\phi_t)\|^2 \\ &\leq \frac{1 - \alpha h_{\min}}{2\alpha} d^2(\phi_t, \bar{\phi}_t) - \frac{1}{2\alpha} d^2(\mathbf{w}_t, \bar{\mathbf{w}}_t) + \frac{\zeta\alpha}{2} \|\nabla P(\phi_t)\|^2, \end{aligned} \quad (57)$$

where the second inequality is from Corollary 5.4, for the update  $\mathbf{w}_t = \exp_{\phi_t}(-\alpha \nabla P(\phi_t))$  from (7) and the third inequality uses the fact  $d^2(\mathbf{w}_t, \bar{\mathbf{w}}_t) = \sum_{k=1}^K d^2(\mathbf{w}_{k,t}, \mathbf{w}_{c,t}) \leq \sum_{k=1}^K d^2(\mathbf{w}_{k,t}, \phi_{m,t}) = d^2(\mathbf{w}_t, \bar{\phi}_t)$  where  $\bar{\mathbf{w}}_t = \text{col}\{\mathbf{w}_{m,t}, \dots, \mathbf{w}_{m,t}\}$  with  $\mathbf{w}_{m,t}$  being the Fréchet mean of  $\mathbf{w}_t$ .

From Corollary 5.4, for the update  $\phi_{t+1} = \exp_{\mathbf{w}_t}(-\mu \widehat{\nabla J}(\mathbf{w}_t))$  in (6), we have

$$\langle -\widehat{\nabla J}(\mathbf{w}_t), \exp_{\mathbf{w}_t}^{-1}(\bar{\mathbf{w}}_t) \rangle \leq \frac{1}{2\mu} d^2(\mathbf{w}_t, \bar{\mathbf{w}}_t) - \frac{1}{2\mu} d^2(\phi_{t+1}, \bar{\mathbf{w}}_t) + \frac{\zeta\mu}{2} \|\widehat{\nabla J}(\mathbf{w}_t)\|^2. \quad (58)$$

Take the expectation on the previous result w.r.t.  $\{\mathbf{x}_s\}_{s=0}^t$  and consider (17) in Assumption 5.7, we obtain

$$\begin{aligned} \mathbb{E}\{\langle -\widehat{\nabla J}(\mathbf{w}_t), \exp_{\mathbf{w}_t}^{-1}(\bar{\mathbf{w}}_t) \rangle\} &= \mathbb{E}\{\langle -\mathbb{E}\{\widehat{\nabla J}(\mathbf{w}_t) | \mathcal{F}_t\}, \exp_{\mathbf{w}_t}^{-1}(\bar{\mathbf{w}}_t) \rangle\} \\ &= \mathbb{E}\{\langle -\nabla J(\mathbf{w}_t), \exp_{\mathbf{w}_t}^{-1}(\bar{\mathbf{w}}_t) \rangle\}, \end{aligned} \quad (59)$$

Combining (58) and (59), we can write

$$\mathbb{E}\{\langle -\nabla J(\mathbf{w}_t), \exp_{\mathbf{w}_t}^{-1}(\bar{\mathbf{w}}_t) \rangle\} \leq \frac{1}{2\mu} \mathbb{E} d^2(\mathbf{w}_t, \bar{\mathbf{w}}_t) - \frac{1}{2\mu} \mathbb{E} d^2(\phi_{t+1}, \bar{\mathbf{w}}_t) + \frac{\zeta\mu}{2} \mathbb{E} \|\widehat{\nabla J}(\mathbf{w}_t)\|^2, \quad (60)$$

Consider  $J$  to be a geodesically convex function under Assumption 5.5. Using (13), taking its expectation w.r.t.  $\{\mathbf{x}_s\}_{s=0}^t$  and combining the result with (60), we further write

$$\begin{aligned} \mathbb{E}\{J(\mathbf{w}_t) - J(\bar{\mathbf{w}}_t)\} &\leq \mathbb{E}\{\langle -\nabla J(\mathbf{w}_t), \exp_{\mathbf{w}_t}^{-1}(\bar{\mathbf{w}}_t) \rangle\} \\ &\leq \frac{1}{2\mu} \mathbb{E} d^2(\mathbf{w}_t, \bar{\mathbf{w}}_t) - \frac{1}{2\mu} \mathbb{E} d^2(\phi_{t+1}, \bar{\mathbf{w}}_t) + \frac{\zeta\mu}{2} \mathbb{E} \|\widehat{\nabla J}(\mathbf{w}_t)\|^2. \end{aligned} \quad (61)$$

Using  $J(\bar{\mathbf{w}}_t) \leq J(\mathbf{w}_t)$  in Lemma 5.9, from (61), we have

$$\begin{aligned} -\mathbb{E} d^2(\mathbf{w}_t, \bar{\mathbf{w}}_t) &\leq -\mathbb{E} d^2(\phi_{t+1}, \bar{\phi}_{t+1}) + \zeta\mu^2 \mathbb{E} \|\widehat{\nabla J}(\mathbf{w}_t)\|^2 \\ &\leq -\mathbb{E} d^2(\phi_{t+1}, \bar{\phi}_{t+1}) + 2\zeta\mu^2 \mathbb{E} \|\nabla J(\mathbf{w}_t)\|^2 + 2\zeta\mu^2 \mathbb{E} \{\|\widehat{\nabla J}(\mathbf{w}_t) - \nabla J(\mathbf{w}_t)\|^2\} \\ &= -\mathbb{E} d^2(\phi_{t+1}, \bar{\phi}_{t+1}) + 2\zeta\mu^2 \mathbb{E} \|\nabla J(\mathbf{w}_t)\|^2 + 2\zeta\mu^2 \mathbb{E} \{\mathbb{E}\{\|\widehat{\nabla J}(\mathbf{w}_t) - \nabla J(\mathbf{w}_t)\|^2 | \mathcal{F}_t\}\} \\ &\leq -\mathbb{E} d^2(\phi_{t+1}, \bar{\phi}_{t+1}) + 2\zeta\mu^2 \mathbb{E} \|\nabla J(\mathbf{w}_t)\|^2 + 2\zeta\mu^2 \sigma^2, \end{aligned} \quad (62)$$

where we use the fact  $\frac{1}{2}(a+b)^2 \leq a^2 + b^2$  in the second equality, and (18) from Assumption 5.7 in the third inequality. Take expectation of (57) w.r.t.  $\{\mathbf{x}_s\}_{s=0}^t$  and combine the result with (62), we obtain

$$\mathbb{E} P(\phi_t) \leq \frac{1 - \alpha h_{\min}}{2\alpha} \mathbb{E} d^2(\phi_t, \bar{\phi}_t) - \frac{1}{2\alpha} \mathbb{E} d^2(\phi_{t+1}, \bar{\phi}_{t+1}) + \frac{\zeta\alpha}{2} \mathbb{E} \|\nabla P(\phi_t)\|^2 + \frac{\zeta\mu^2}{\alpha} \mathbb{E} \|\nabla J(\mathbf{w}_t)\|^2 + \frac{\zeta\mu^2 \sigma^2}{\alpha}. \quad (63)$$



Multiplying (22) in Lemma 5.11 by  $2\zeta$  and summing the result to (63), and considering the upper bound of  $\|\nabla J(\mathbf{w}_t)\|^2$  given in Lemma 5.10, we have

$$2\zeta \mathbb{E}P(\phi_{t+1}) - (2\zeta - 1)\mathbb{E}P(\phi_t) \leq \frac{1 - \alpha h_{\min}}{2\alpha} \mathbb{E}d^2(\phi_t, \bar{\phi}_t) - \frac{1}{2\alpha} \mathbb{E}d^2(\phi_{t+1}, \bar{\phi}_{t+1}) + \frac{11\zeta\mu^2}{\alpha} G^2 + \frac{3\zeta\mu^2}{\alpha} \sigma^2. \quad (64)$$

Multiplying (64) by  $(1 - \tau)^{-t}$ , we have:

$$\begin{aligned} (1 - \tau)^{-t} 2\zeta \mathbb{E}P(\phi_{t+1}) - (1 - \tau)^{-t} (1 - \frac{1}{2\zeta}) 2\zeta \mathbb{E}P(\phi_t) &\leq (1 - \tau)^{-t} \frac{1 - \alpha h_{\min}}{2\alpha} \mathbb{E}d^2(\phi_t, \bar{\phi}_t) \\ &\quad - (1 - \tau)^{-t} \frac{1}{2\alpha} \mathbb{E}d^2(\phi_{t+1}, \bar{\phi}_{t+1}) \\ &\quad + (1 - \tau)^{-t} \frac{11\zeta\mu^2}{\alpha} G^2 + (1 - \tau)^{-t} \frac{3\zeta\mu^2}{\alpha} \sigma^2. \end{aligned} \quad (65)$$

Now we sum (65) from  $t = 0$  to  $t = s - 1$ . To simplify the summation, we consider the case  $t = 0$  and  $t \geq 1$  separately as we can get a simpler upper bound in the latter case. Consider the case  $t = 0$ , which is simple. From (64) we have:

$$2\zeta \mathbb{E}P(\phi_1) - (2\zeta - 1)\mathbb{E}P(\phi_0) \leq \frac{1 - \alpha h_{\min}}{2\alpha} \mathbb{E}d^2(\phi_0, \bar{\phi}_0) - \frac{1}{2\alpha} \mathbb{E}d^2(\phi_1, \bar{\phi}_1) + \frac{11\zeta\mu^2}{\alpha} G^2 + \frac{3\zeta\mu^2}{\alpha} \sigma^2. \quad (66)$$

For the case  $t \geq 1$ , inspired by (Zhang & Sra, 2016), let  $\tau = \min\{\frac{1}{2\zeta}, \alpha h_{\min}\}$ , this implies  $\tau \leq \frac{1}{2\zeta}$  and  $\tau \leq \alpha h_{\min}$ . Consider  $\alpha \leq h_{\max}^{-1} < h_{\min}^{-1}$ , we have  $\tau \in (0, 1)$ . For  $t \geq 1$ , from (65) we can obtain:

$$\begin{aligned} (1 - \tau)^{-t} 2\zeta \mathbb{E}P(\phi_{t+1}) - (1 - \tau)^{-(t-1)} 2\zeta \mathbb{E}P(\phi_t) &\leq (1 - \tau)^{-(t-1)} \frac{1}{2\alpha} \mathbb{E}d^2(\phi_t, \bar{\phi}_t) - (1 - \tau)^{-t} \frac{1}{2\alpha} \mathbb{E}d^2(\phi_{t+1}, \bar{\phi}_{t+1}) \\ &\quad + (1 - \tau)^{-t} \frac{11\zeta\mu^2}{\alpha} G^2 + (1 - \tau)^{-t} \frac{3\zeta\mu^2}{\alpha} \sigma^2. \end{aligned} \quad (67)$$

Finally, summing (65) over  $t$  from  $t = 0$  to  $t = s - 1$ , and using the previous results, we have:

$$\begin{aligned} (1 - \tau)^{-(s-1)} 2\zeta \mathbb{E}P(\phi_s) - (2\zeta - 1)\mathbb{E}P(\phi_0) &\leq \frac{1 - \alpha h_{\min}}{2\alpha} \mathbb{E}d^2(\phi_0, \bar{\phi}_0) - (1 - \tau)^{-(s-1)} \frac{1}{2\alpha} \mathbb{E}d^2(\phi_s, \bar{\phi}_s) \\ &\quad + \sum_{t=0}^{s-1} (1 - \tau)^{-t} \frac{11\zeta\mu^2}{\alpha} G^2 + \sum_{t=0}^{s-1} (1 - \tau)^{-t} \frac{3\zeta\mu^2}{\alpha} \sigma^2 \\ &\leq \frac{D^2}{2\alpha} + \sum_{t=0}^{s-1} (1 - \tau)^{-t} \frac{11\zeta\mu^2}{\alpha} G^2 + \sum_{t=0}^{s-1} (1 - \tau)^{-t} \frac{3\zeta\mu^2}{\alpha} \sigma^2, \end{aligned} \quad (68)$$

where the second inequality drops the negative terms and plugs in  $d(\phi_0, \bar{\phi}_0) \leq D$  (Assumption 5.1).

Define  $\gamma_4(\beta) \triangleq \exp_{\phi_0}(\beta \exp_{\phi_0}^{-1}(\bar{\phi}_0))$  as the minimal geodesic from  $\phi_0$  to  $\bar{\phi}_0$ . Using the second-order Taylor expansion of  $\beta \mapsto P(\gamma_4(\beta))$  around  $\beta = 1$ , considering  $P(\bar{\phi}_0) = 0$ ,  $\nabla P(\bar{\phi}_0) = 0$ , we have (Afsari et al., 2013):

$$\begin{aligned} P(\phi_0) &\leq P(\bar{\phi}_0) + \langle \nabla P(\bar{\phi}_0), \exp_{\phi_0}^{-1}(\bar{\phi}_0) \rangle + \frac{h_{\max} \|\exp_{\phi_0}^{-1}(\bar{\phi}_0)\|^2}{2} \\ &= \frac{h_{\max}}{2} d^2(\phi_0, \bar{\phi}_0). \end{aligned} \quad (69)$$

This ensures  $P(\phi_0) \leq \frac{h_{\max}}{2} D^2 \leq \frac{D^2}{2\alpha}$  since  $d(\phi_0, \bar{\phi}_0) \leq D$  and  $\alpha \in (0, h_{\max}^{-1}]$ , one can thus obtain from (68) that

$$\begin{aligned} \mathbb{E}P(\phi_s) &\leq \frac{(1 - \tau)^{-(s-1)} D^2}{2\alpha} + \sum_{t=0}^{s-1} (1 - \tau)^{-t} \frac{11\mu^2}{2\alpha} G^2 + \sum_{t=0}^{s-1} (1 - \tau)^{-t} \frac{3\mu^2}{2\alpha} \sigma^2 \\ &\leq \frac{(1 - \tau)^{-(s-1)} D^2}{2\alpha} + \sum_{t=0}^{\infty} (1 - \tau)^{-t} \frac{11\mu^2}{2\alpha} G^2 + \sum_{t=0}^{\infty} (1 - \tau)^{-t} \frac{3\mu^2}{2\alpha} \sigma^2 \end{aligned}$$

$$\begin{aligned}
 &\leq \frac{(1-\tau)^{(s-1)}D^2}{2\alpha} + \frac{11\mu^2}{2\alpha\tau}G^2 + \frac{3\mu^2}{2\alpha\tau}\sigma^2 \\
 &\leq \frac{11\mu^2}{2\alpha\tau}G^2 + \frac{3\mu^2}{\alpha\tau}\sigma^2,
 \end{aligned} \tag{70}$$

where the last inequality holds whenever:

$$\begin{aligned}
 \frac{(1-\tau)^{(s-1)}D^2}{2\alpha} &\leq \frac{3\mu^2}{2\alpha\tau}\sigma^2 \iff (1-\tau)^{(s-1)} \leq \frac{3\mu^2}{\tau D^2}\sigma^2 \\
 &\iff (s-1)\log(1-\tau) \leq 2\log(\mu) + O(1) \\
 &\iff s \leq \frac{2\log(\mu)}{\log(1-\tau)} + O(1).
 \end{aligned} \tag{71}$$

We conclude that

$$\mathbb{E}\{P(\phi_s)\} \leq \frac{11\mu^2}{2\alpha\tau}G^2 + \frac{3\mu^2}{\alpha\tau}\sigma^2, \tag{72}$$

with sufficiently small step sizes  $\mu$  after sufficient iterations  $s_o$ , where

$$s_o = \frac{2\log(\mu)}{\log(1-\tau)} + O(1) = O(\mu^{-1}) \tag{73}$$

where the second equality follows since  $\lim_{\mu \rightarrow 0} \mu \log(\mu) = 0$ , which means that the magnitude of  $\log(\mu)$  can be bounded above by a constant multiple of  $\mu^{-1}$  for  $\mu \rightarrow 0$ .

## B.2. Theorem 5.15

Denote  $\Delta_t = J(\mathbf{w}_t) - J(\mathbf{w}^*)$ , from Lemma 5.14, we have:

$$\mathbb{E}\Delta_{t+1} - \mathbb{E}\Delta_t \leq -\frac{\mu}{4}\mathbb{E}\|\widehat{\nabla J}(\mathbf{w}_t)\|^2 + \frac{5\alpha^2}{\mu}\mathbb{E}\|\nabla P(\phi_{t+1})\|^2. \tag{74}$$

From Corollary 5.4, for the update  $\phi_{t+1} = \exp_{\mathbf{w}_t}(-\mu\widehat{\nabla J}(\mathbf{w}_t))$  in (6), we have

$$\langle -\widehat{\nabla J}(\mathbf{w}_t), \exp_{\mathbf{w}_t}^{-1}(\mathbf{w}^*) \rangle \leq \frac{1}{2\mu}d^2(\mathbf{w}_t, \mathbf{w}^*) - \frac{1}{2\mu}d^2(\phi_{t+1}, \mathbf{w}^*) + \frac{\zeta\mu}{2}\|\widehat{\nabla J}(\mathbf{w}_t)\|^2. \tag{75}$$

Take the expectation on the previous result w.r.t.  $\{\mathbf{x}_s\}_{s=0}^t$  and consider (17) in Assumption 5.7, we obtain

$$\begin{aligned}
 \mathbb{E}\{\langle -\widehat{\nabla J}(\mathbf{w}_t), \exp_{\mathbf{w}_t}^{-1}(\mathbf{w}^*) \rangle\} &= \mathbb{E}\{\langle -\mathbb{E}\{\widehat{\nabla J}(\mathbf{w}_t)|\mathcal{F}_t\}, \exp_{\mathbf{w}_t}^{-1}(\mathbf{w}^*) \rangle\} \\
 &= \mathbb{E}\{\langle -\nabla J(\mathbf{w}_t), \exp_{\mathbf{w}_t}^{-1}(\mathbf{w}^*) \rangle\},
 \end{aligned} \tag{76}$$

Combining (75) and (76), we have

$$\mathbb{E}\{\langle -\nabla J(\mathbf{w}_t), \exp_{\mathbf{w}_t}^{-1}(\mathbf{w}^*) \rangle\} \leq \frac{1}{2\mu}\mathbb{E}d^2(\mathbf{w}_t, \mathbf{w}^*) - \frac{1}{2\mu}\mathbb{E}d^2(\phi_{t+1}, \mathbf{w}^*) + \frac{\zeta\mu}{2}\mathbb{E}\|\widehat{\nabla J}(\mathbf{w}_t)\|^2, \tag{77}$$

Consider that  $J$  is a geodesically convex function under Assumption 5.5, from (77) one can obtain

$$\begin{aligned}
 \mathbb{E}\Delta_t &= \mathbb{E}\{J(\mathbf{w}_t) - J(\mathbf{w}^*)\} \leq \mathbb{E}\{\langle -\nabla J(\mathbf{w}_t), \exp_{\mathbf{w}_t}^{-1}(\mathbf{w}^*) \rangle\} \\
 &\leq \frac{1}{2\mu}\mathbb{E}d^2(\mathbf{w}_t, \mathbf{w}^*) - \frac{1}{2\mu}\mathbb{E}d^2(\phi_{t+1}, \mathbf{w}^*) + \frac{\zeta\mu}{2}\mathbb{E}\|\widehat{\nabla J}(\mathbf{w}_t)\|^2.
 \end{aligned} \tag{78}$$

Now we need to upper bound  $-d_{\mathcal{M}}^2(\phi_{t+1}, \mathbf{w}^*)$ . From Corollary 5.4, for the update  $\mathbf{w}_{t+1} = \exp_{\phi_{t+1}}(-\alpha\nabla P(\phi_{t+1}))$  in (7), we have

$$d^2(\mathbf{w}_{t+1}, \mathbf{w}^*) - d^2(\phi_{t+1}, \mathbf{w}^*) \leq \zeta\alpha^2\|\nabla P(\phi_{t+1})\|^2 + 2\alpha\langle \nabla P(\phi_{t+1}), \exp_{\phi_{t+1}}^{-1}(\mathbf{w}^*) \rangle$$

$$\begin{aligned} &\leq \zeta \alpha^2 \|\nabla P(\phi_{t+1})\|^2 + 2\alpha(P(\mathbf{w}^*) - P(\phi_{t+1})) \\ &\leq \zeta \alpha^2 \|\nabla P(\phi_{t+1})\|^2, \end{aligned} \quad (79)$$

where the second inequality uses the convexity property of  $P$  with (13), and the third inequality uses the fact  $P(\mathbf{w}^*) = 0$  and  $P(\phi_{t+1}) \geq 0$ .

Take the expectation on the previous result w.r.t.  $\{\mathbf{x}_s\}_{s=0}^t$ , and combine the result with (78), we have

$$\mathbb{E}\Delta_t \leq \frac{1}{2\mu} \mathbb{E}d^2(\mathbf{w}_t, \mathbf{w}^*) - \frac{1}{2\mu} \mathbb{E}d^2(\mathbf{w}_{t+1}, \mathbf{w}^*) + \frac{\zeta\mu}{2} \mathbb{E}\|\widehat{\nabla J}(\mathbf{w}_t)\|^2 + \frac{\zeta\alpha^2}{2\mu} \mathbb{E}\|\nabla P(\phi_{t+1})\|^2. \quad (80)$$

Multiplying (74) by  $2\zeta$  and adding to (80), we have:

$$2\zeta \mathbb{E}\Delta_{t+1} - (2\zeta - 1)\mathbb{E}\Delta_t \leq \frac{1}{2\mu} \mathbb{E}d^2(\mathbf{w}_t, \mathbf{w}^*) - \frac{1}{2\mu} \mathbb{E}d^2(\mathbf{w}_{t+1}, \mathbf{w}^*) + \frac{21\zeta\alpha^2}{2\mu} \mathbb{E}\|\nabla P(\phi_{t+1})\|^2. \quad (81)$$

From Corollary 5.13, we know  $\mathbb{E}\|\nabla P(\phi_t)\|^2$  converges to a small value after sufficient iterations  $s_o$ , now we tend to study the convergence of our algorithm after  $s_o$  iterations. For  $t \geq s_o$  where  $s_o$  defined in (24), from (27) one can rewrite (81) as

$$2\zeta \mathbb{E}\Delta_{t+1} - (2\zeta - 1)\mathbb{E}\Delta_t \leq \frac{1}{2\mu} \mathbb{E}d^2(\mathbf{w}_t, \mathbf{w}^*) - \frac{1}{2\mu} \mathbb{E}d^2(\mathbf{w}_{t+1}, \mathbf{w}^*) + \frac{231\zeta\alpha\mu}{2\tau} G^2 + \frac{63\zeta\alpha\mu}{\tau} \sigma^2. \quad (82)$$

Summing (82) from  $t = s_o$  to  $t = s - 1$  and plugging in  $d(\mathbf{w}_{s_o}, \mathbf{w}^*) \leq D$ , we obtain

$$\begin{aligned} 2\zeta \mathbb{E}\Delta_s + \sum_{t=s_o+1}^{s-1} \mathbb{E}\Delta_t &\leq (2\zeta - 1)\mathbb{E}\Delta_{s_o} + \frac{1}{2\mu} d^2(\mathbf{w}_{s_o}, \mathbf{w}^*) - \frac{1}{2\mu} \mathbb{E}d^2(\mathbf{w}_s, \mathbf{w}^*) + (s - s_o) \left( \frac{231\zeta\alpha\mu}{2\tau} G^2 + \frac{63\zeta\alpha\mu}{\tau} \sigma^2 \right) \\ &\leq (2\zeta - 1)\mathbb{E}\Delta_{s_o} + \frac{1}{2\mu} d^2(\mathbf{w}_{s_o}, \mathbf{w}^*) + (s - s_o) \left( \frac{231\zeta\alpha\mu}{2\tau} G^2 + \frac{63\zeta\alpha\mu}{\tau} \sigma^2 \right) \\ &\leq (2\zeta - 1)\mathbb{E}\Delta_{s_o} + \frac{D^2}{2\mu} + (s - s_o) \left( \frac{231\zeta\alpha\mu}{2\tau} G^2 + \frac{63\zeta\alpha\mu}{\tau} \sigma^2 \right). \end{aligned} \quad (83)$$

Recall the geodesic  $L$ -smoothness of  $J$  in Assumption 5.6 and plugging into  $d(\mathbf{w}_{s_o}, \mathbf{w}^*) \leq D$  and  $\nabla J(\mathbf{w}^*) = 0$ , we have:

$$\begin{aligned} \Delta_{s_o} &= J(\mathbf{w}_{s_o}) - J(\mathbf{w}^*) \leq \langle \nabla J(\mathbf{w}^*), \exp_{\mathbf{w}^*}^{-1}(\mathbf{w}_{s_o}) \rangle + \frac{L}{2} \|\exp_{\mathbf{w}^*}^{-1}(\mathbf{w}_{s_o})\|^2 \\ &= \frac{L}{2} d^2(\mathbf{w}_{s_o}, \mathbf{w}^*) \leq \frac{LD^2}{2}. \end{aligned} \quad (84)$$

This ensures  $\Delta_{s_o} \leq \frac{LD^2}{2} \leq \frac{D^2}{2\mu}$  since  $\mu \leq L^{-1}$ , so that from (83) one can obtain

$$2\zeta \mathbb{E}\Delta_s + \sum_{t=s_o+1}^{s-1} \mathbb{E}\Delta_t \leq \zeta LD^2 + (s - s_o) \left( \frac{231\zeta\alpha\mu}{2\tau} G^2 + \frac{63\zeta\alpha\mu}{\tau} \sigma^2 \right). \quad (85)$$

Here the term  $\Delta_s$  does not cancel nicely due to the presence of the curvature term  $\zeta$ , which necessitates the use of a Lyapunov function as in (Zhang & Sra, 2016). Introduce auxiliary variables  $\mathbf{w}'_{s_o+1} = \mathbf{w}_{s_o+1}$  and  $\mathbf{w}'_{t+1} = \exp_{\mathbf{w}'_t} \left( \frac{1}{t-s_o+1} \exp_{\mathbf{w}'_t}^{-1}(\mathbf{w}_{t+1}) \right)$  for  $s_o + 1 \leq t \leq s - 2$ ,  $\mathbf{w}'_{s+1} = \exp_{\mathbf{w}'_{s-1}} \left( \frac{2\zeta}{2\zeta+s-s_o-1} \exp_{\mathbf{w}'_{s-1}}^{-1}(\mathbf{w}_s) \right)$ , repeatedly consider (12) in Assumption 5.5 (geodesic convexity of  $J$ ), for  $s \geq s_o + 1$ , we have

$$\begin{aligned} J(\mathbf{w}'_{s-1}) &\leq \frac{s-s_o-2}{s-s_o-1} J(\mathbf{w}'_{s-2}) + \frac{1}{s-s_o-1} J(\mathbf{w}_{s-1}) \\ &\leq \frac{s-s_o-2}{s-s_o-1} \left( \frac{s-s_o-3}{s-s_o-2} J(\mathbf{w}'_{s-3}) + \frac{1}{s-s_o-2} J(\mathbf{w}_{s-2}) \right) + \frac{1}{s-s_o-1} J(\mathbf{w}_{s-1}) \\ &\leq \dots \leq \frac{1}{s-s_o-1} \sum_{t=s_o+1}^{s-1} J(\mathbf{w}_t). \end{aligned} \quad (86)$$

Denote  $\Delta'_s = J(\mathbf{w}'_s) - J(\mathbf{w}^*)$ , we have  $\mathbb{E}\Delta'_{s-1} \leq \frac{1}{s-s_o-1} \sum_{t=s_o+1}^{s-1} \mathbb{E}\Delta_t$ . Again, consider the geodesic convexity of  $J$  in (12) of Assumption 5.5, and we can further write

$$\begin{aligned} \mathbb{E}\Delta'_s &= \mathbb{E}\{J(\mathbf{w}'_s) - J(\mathbf{w}^*)\} \leq \mathbb{E}\left\{\frac{s-s_o-1}{2\zeta+s-s_o-1}J(\mathbf{w}'_{s-1}) + \frac{2\zeta}{2\zeta+s-s_o-1}J(\mathbf{w}_s) - J(\mathbf{w}^*)\right\} \\ &= \frac{2\zeta\mathbb{E}\Delta_s + (s-s_o-1)\mathbb{E}\Delta'_{s-1}}{2\zeta+s-s_o-1} \\ &\leq \frac{2\zeta\mathbb{E}\Delta_s + \sum_{t=s_o+1}^{s-1} \mathbb{E}\Delta_t}{2\zeta+s-s_o-1}. \end{aligned} \quad (87)$$

Plug the upper bound of  $2\zeta\mathbb{E}\Delta_s + \sum_{t=s_o+1}^{s-1} \mathbb{E}\Delta_t$  in (83) into the above result, for  $s \geq s_o + 1$ , we have

$$\mathbb{E}\Delta'_s \leq \frac{\zeta LD^2}{2\zeta+s-s_o-1} + \frac{s-s_o}{2\zeta+s-s_o-1} \left( \frac{231\zeta\alpha\mu}{2\tau} G^2 + \frac{63\zeta\alpha\mu}{\tau} \sigma^2 \right). \quad (88)$$

## C. Examples of Riemannian manifolds

### C.1. Grassmann manifold

The Grassmann manifold  $\mathcal{G}_n^p$ , a set of  $p$ -dimensional linear subspaces of  $\mathbb{R}^n$ , can be regarded as a smooth quotient manifold of the Stiefel manifold  $\mathcal{S}_n^p = \{\mathbf{U} \in \mathbb{R}^{n \times p} : \mathbf{U}^T \mathbf{U} = \mathbf{I}_p\}$ , i.e.,  $\mathcal{G}_n^p = \mathcal{S}_n^p / \mathcal{O}_p = \{\pi(\mathbf{U}) : \mathbf{U} \in \mathcal{S}_n^p\}$  where  $\mathcal{O}_p = \{\mathbf{U} \in \mathbb{R}^{p \times p} : \mathbf{U}^T \mathbf{U} = \mathbf{I}_p\}$  is the orthogonal group and  $\pi : \mathcal{S}_n^p \rightarrow \mathcal{G}_n^p$  is the map  $\pi(\mathbf{U}) = \{\mathbf{U}\mathbf{O} : \mathbf{O} \in \mathcal{O}_p\}$ . The geodesic distance between two subspaces  $\pi(\mathbf{U}_1)$  and  $\pi(\mathbf{U}_2)$  of  $\mathcal{G}_n^p$ , spanned by orthonormal matrices  $\mathbf{U}_1$  and  $\mathbf{U}_2$ , is defined as follows (Edelman et al., 1998):

$$d_{\mathcal{G}_n^p}(\mathbf{U}_1, \mathbf{U}_2) = \|\cos^{-1}(\boldsymbol{\theta})\|_2, \quad (89)$$

where  $\boldsymbol{\theta} \in \mathbb{R}^p$  contains the singular values of  $\mathbf{U}_1^T \mathbf{U}_2$ , namely, it is related to its singular value decomposition (SVD) as  $\mathbf{U}_1^T \mathbf{U}_2 = \mathbf{V}_1^T \text{diag}(\boldsymbol{\theta}) \mathbf{V}_2$ . Define  $\bar{f} : \mathcal{S}_n^p \rightarrow \mathbb{R}$ , we have  $f(\pi(\mathbf{U})) = \bar{f}(\mathbf{U})$  for all  $\pi(\mathbf{U}) \in \mathcal{G}_n^p$ . The Riemannian gradient  $\nabla f$  at  $\pi(\mathbf{U}) \in \mathcal{G}_n^p$  is given by:

$$\nabla f(\pi(\mathbf{U})) = \nabla \bar{f}(\mathbf{U}) = \mathbf{P}_{\mathbf{U}}^{\mathcal{G}_n^p}(\mathbf{G}), \quad (90)$$

with  $\mathbf{P}_{\mathbf{U}}^{\mathcal{G}_n^p}(\mathbf{G}) = (\mathbf{I} - \mathbf{U}\mathbf{U}^T)\mathbf{G}$ , where  $\mathbf{G} \in \mathbb{R}^{n \times p}$  is the Euclidean gradient of  $\bar{f}$  at  $\mathbf{U}$ . Let  $\boldsymbol{\xi} \in T_{\pi(\mathbf{U})}\mathcal{G}_n^p$ , and let  $\mathbf{X}\boldsymbol{\Sigma}\mathbf{Y} = \mathbf{U} + \boldsymbol{\xi}$  be the thin SVD of  $\mathbf{U} + \boldsymbol{\xi} \in \mathbb{R}^{n \times p}$ . A numerically stable retraction  $R_{\pi(\mathbf{U})} : T_{\pi(\mathbf{U})}\mathcal{G}_n^p \rightarrow \mathcal{G}_n^p$  on  $\mathcal{G}_n^p$  is given by (Boumal, 2023):

$$R_{\pi(\mathbf{U})}(\boldsymbol{\xi}) = \pi(\mathbf{X}\mathbf{Y}^T). \quad (91)$$

### C.2. The manifold of SPD matrices

The geodesic distance of  $\mathcal{S}_n^{++}$  between two SPD matrices  $\boldsymbol{\Sigma}_1$  and  $\boldsymbol{\Sigma}_2 \in \mathcal{S}_n^{++}$  can be computed in closed form (Pennec et al., 2006) as:

$$d_{\mathcal{S}_n^{++}}(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2) = \left\| \log(\boldsymbol{\Sigma}_2^{-\frac{1}{2}} \boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}_2^{-\frac{1}{2}}) \right\|_F, \quad (92)$$

where  $\|\cdot\|_F$  denotes the Frobenius norm. The Riemannian gradient  $\nabla f$  at  $\boldsymbol{\Sigma} \in \mathcal{S}_n^{++}$  is given by:

$$\nabla f(\boldsymbol{\Sigma}) = \boldsymbol{\Sigma} \text{sym}(\mathbf{G}) \boldsymbol{\Sigma}, \quad (93)$$

with  $\mathbf{G} \in \mathbb{R}^{p \times p}$  the Euclidean gradient of function  $f$  at  $\boldsymbol{\Sigma}$  and  $\text{sym}(\mathbf{G}) = \frac{1}{2}(\mathbf{G}^T + \mathbf{G})$ . In practice, the Euclidean gradient can be easily computed using automatic differentiation tools. Let  $\boldsymbol{\xi} \in T_{\boldsymbol{\Sigma}}\mathcal{S}_n^{++}$ . A retraction  $R_{\boldsymbol{\Sigma}, \mathcal{S}_n^{++}} : T_{\boldsymbol{\Sigma}}\mathcal{S}_n^{++} \rightarrow \mathcal{S}_n^{++}$  is defined as:

$$R_{\boldsymbol{\Sigma}, \mathcal{S}_n^{++}}(\boldsymbol{\xi}) = \boldsymbol{\Sigma} + \boldsymbol{\xi} + \frac{1}{2}\boldsymbol{\xi}\boldsymbol{\Sigma}^{-1}\boldsymbol{\xi}. \quad (94)$$

This retraction is a second-order approximation of the exponential mapping.



## D. Additional experimental results

### D.1. Inefficiency of the method (Wang et al., 2024b)

In Section 2, we argue that the algorithm in (Wang et al., 2024b) is inefficient due to the inner-loop optimization when minimizing the penalty term  $P(\phi_t)$ . To support this claim, we compare the MSD performance and runtime between the work in (Wang et al., 2024b) (denoted as “Inefficient Riemannian diffusion”) and the proposed algorithm. We examine these two algorithms for distributed PCA on synthetic data in the same setting as in Subsection 7.1, and produce the results as in Figure 6. From these results, we can see that while the performance of these two algorithms is nearly identical, the proposed algorithm achieves a significantly reduced runtime. These experiments were performed on a computer with an Apple M4 Pro processor and 24GB of RAM.

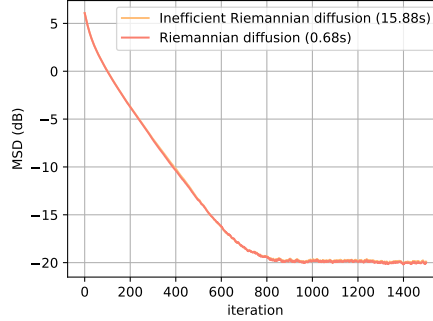


Figure 6. Illustration of MSD performance and runtime per Monte Carlo run of the (inefficient) Riemannian diffusion adaptation algorithms for distributed PCA on synthetic data.

### D.2. Applicability to more networks

To illustrate the applicability of the proposed algorithm to more networks, we randomly generate another graph topology as shown in Figure 7 (left) and select weights with a uniform rule. We test all compared algorithms for both distributed PCA and GMM inference on synthetic data in the same setting as in Section 7 and produce experimental results in Figure 7 (middle and right). From these results, we find the performance of the compared algorithms remains similar to that shown in Figure 2 and Figure 4, which are obtained with the network illustrated in Figure 1.

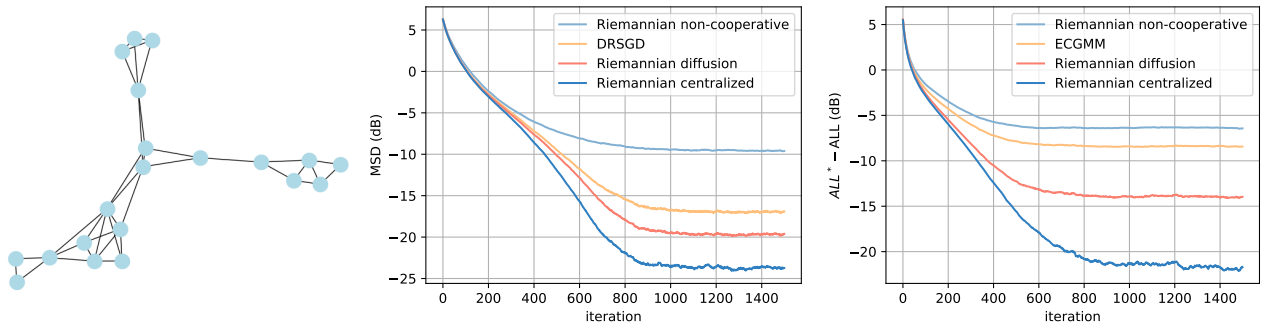


Figure 7. Performance illustration of the compared methods on another network with uniformly distributed weights: graph topology (left); MSD performance for distributed PCA (middle); ALL differences for GMM inference (right) on synthetic data.

### D.3. Impact of step sizes

For the proposed algorithm, the choice of step sizes is critical to control the tradeoff between convergence speed and steady-state performance. We examine the behavior of the proposed algorithm for distributed PCA on synthetic data in the same setting as in Subsection 7.1, and produce the results with different choices of step sizes as shown in Figure 8. It can be observed that larger step sizes tend to accelerate convergence but result in worse performance at steady state.

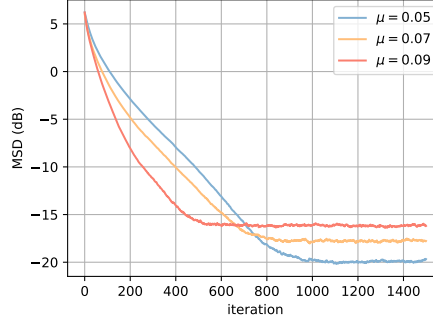


Figure 8. Illustration of MSD performance of the proposed method with different step sizes for distributed PCA on synthetic data.

## E. Computational complexity

The computational complexity of the proposed algorithm on each agent  $k$  involves two contributing terms. The first is the cost of a local adaptation step (3) (i.e., Riemannian SGD on  $J_k$ ), which is denoted by  $T_J$ . The second is the cost of the combination step (4), which involves a gradient step over the loss function  $P_k$  that scales linearly with the number  $N_{\text{neigh},k}$  of neighbors connected to node  $k$  in the graph (that is, with the number of nonzero elements in the coefficients  $c_{k\ell}$ ), which we represent as  $N_{\text{neigh},k} \cdot T_P$ , where  $T_P$  is the cost of computing the inverse of the exponential mapping.  $N_{\text{neigh},k}$  is also known as the *degree* of the vertex  $k$  in the graph  $\mathcal{G}$ . Thus, for each agent  $k$ , we obtain a complexity of  $T_J + N_{\text{neigh},k} \cdot T_P$ . Compared to a non-cooperative setting, we have an overhead cost of  $N_{\text{neigh},k} \cdot T_P$ , which is a function of both on  $T_P$  (which depends on the manifold) and on the number of neighbors connected to node  $k$  (which depends on the graph topology).

This allows us to understand how the complexity scales with the number of agents  $K$ . In the case where the number of neighbors to each node (i.e., their degree in the graph) is constant, the complexity does not increase with  $K$ . On the other hand, in the worst case scenario of a fully connected graph (where each vertex has degree  $K - 1$ , being connected to all other vertices), then the complexity scales linearly with  $K$ , with a coefficient equal to  $T_P$ .

## F. Discussion on the limitations

Our work has two main limitations, which are discussed in the following.

- The theoretical analysis is based on the exponential mapping  $\exp_x$  as in many works in Riemannian optimization, e.g., (Zhang & Sra, 2016), while in practice, a retraction  $R_x$  is used for more efficient computations. A key result in (Bonnabel, 2013) states that  $d(R_x(\mu \cdot v), \exp_x(\mu \cdot v)) = O(\mu^2)$ , meaning that for small  $\mu$ , a retraction closely approximates the exponential map. The main approach to proving convergence with retractions involves showing that the iterates of the algorithm remain close to those of an equivalent version using the exponential map, which holds as  $\mu \rightarrow 0$  (Bonnabel, 2013). This argument typically relies on diminishing step sizes, whereas our analysis is designed for constant step sizes, which are crucial for continuous adaptation and learning. Some works also employ the *pullback* operator  $f \circ R_x$ , i.e., the composition of the cost function  $f$  and a retraction  $R_x$ , to establish convergence. However, these approaches require assumptions that may be less natural, such as the convexity and smoothness of the pullback operator, see Chapter 4 of (Boumal, 2023). Thus, we believe that extending the proposed theoretical analysis based on a retraction is an exciting, though non-trivial, research direction.
- Manifolds without closed-form expressions for retractions, or for the Riemannian gradient, pose challenges to the implementation of the proposed algorithm, as such operations have to be approximated numerically in some way. However, we highlight that this limitation also applies to most existing Riemannian optimization algorithms and is not specific to our work.