

Riemannian Bilevel Optimization

Jiaxiang Li

*Department of Electrical and Computer Engineering
University of Minnesota, Twin Cities
Minneapolis, MN 55455, USA*

LI003755@UMN.EDU

Shiqian Ma

*Department of Computational Applied Math and Operations Research
Rice University
Houston, TX 77005, USA*

SQMA@RICE.EDU

Editor: Lam Nguyen

Abstract

In this work, we consider the bilevel optimization problem on Riemannian manifolds. We inspect the calculation of the hypergradient of such problems on general manifolds and thus enable the utilization of gradient-based algorithms to solve such problems. The calculation of the hypergradient requires utilizing the notion of Riemannian cross-derivative and we inspect the properties and the numerical calculations of Riemannian cross-derivatives. Algorithms in both deterministic and stochastic settings, named respectively RieBO and RieSBO, are proposed that include the existing Euclidean bilevel optimization algorithms as special cases. Numerical experiments on robust optimization on Riemannian manifolds are presented to show the applicability and efficiency of the proposed methods.

Keywords: Riemannian optimization, stochastic optimization, bilevel optimization

1. Introduction

Bilevel optimization has drawn attentions from various fields in optimization and machine learning communities, due to its wide range of applications including meta learning (Rajeswaran et al., 2019; Ji et al., 2020), hyperparameter optimization (Okuno et al., 2021; Yu and Zhu, 2020), reinforcement learning (Konda and Tsitsiklis, 1999; Hong et al., 2023) and signal processing (Kunapuli et al., 2008; Flamary et al., 2014). In this work, we focus on the manifold-constrained bilevel optimization problem, which can be formulated as:

$$\begin{aligned} \min_{x \in \mathcal{M}} \Phi(x) &:= f(x, y^*(x)) \\ \text{s.t. } y^*(x) &= \operatorname{argmin}_{y \in \mathcal{N}} g(x, y), \end{aligned} \quad (1.1)$$

where \mathcal{M} and \mathcal{N} are m and n -dimensional complete Riemannian manifolds, respectively. We also consider the stochastic bilevel optimization which is in the following form:

$$\begin{aligned} \min_{x \in \mathcal{M}} \Phi(x) &= f(x, y^*(x)) := \mathbb{E}_{\xi} [F(x, y^*(x); \xi)] \\ \text{s.t. } y^*(x) &= \operatorname{argmin}_{y \in \mathcal{N}} g(x, y) := \mathbb{E}_{\zeta} [G(x, y; \zeta)], \end{aligned} \quad (1.2)$$

where ξ and ζ are random variables that usually represent the randomness from the data. Such a framework allows us to utilize the stochastic gradient methods to get a desired convergence result with only a noisy estimate of the gradients for f and g . Here we also assume that g is a (geodesically) strongly convex function with respect to y in both (1.1) and (1.2) so that the solution to the lower level problem $y^*(x)$ is well-defined.

Notice that the original (Euclidean) bilevel optimization is a special case of (1.1) by taking the manifolds as the Euclidean spaces with the same dimensions:

$$\begin{aligned} \min_{x \in \mathbb{R}^m} \Phi(x) &:= f(x, y^*(x)) \\ \text{s.t. } y^*(x) &= \operatorname{argmin}_{y \in \mathbb{R}^n} g(x, y), \end{aligned} \tag{1.3}$$

where f and g are assumed to be continuously differentiable. It is worth noticing that the objective function $\Phi(x)$ is still nonconvex even if we impose convexity assumptions on f , which makes such a problem hard to tackle, let alone the more complicated manifold-constraint problems, namely (1.1) and (1.2).

There has been an extensive study on the Euclidean bilevel optimization (Ji et al., 2021; Hong et al., 2023; Chen et al., 2021; Ghadimi and Wang, 2018). On the algorithmic sense, the bilevel optimization seeks to obtain a first-order ϵ -stationary point (Definition 11 and 19 for deterministic and stochastic cases, respectively) with the access to the gradient oracle of f and g , as well as the Jacobian- and Hessian-vector product, i.e. $\nabla_x \nabla_y g(x, y)v$ and $\nabla_y^2 g(x, y)v$, respectively. To find an ϵ -stationary point, we denote the number of calls to the gradient oracle of f and g as $\text{Gc}(f, \epsilon)$ and $\text{Gc}(g, \epsilon)$, correspondingly; similarly we have the notation JV and HV for the number of oracle calls for the Jacobian- and Hessian-vector product. In the Euclidean setting, We have Tables 1 and 2 as the summary for the oracle calls to achieve an ϵ -stationary point for deterministic and stochastic cases, correspondingly (and we denote κ as the condition number of the lower level strongly convex problem).

Algorithm	BA	AID-BiO	ITD-BiO*
y -update	GD	GD	GD
$\text{Gc}(f, \epsilon)$	$\mathcal{O}(\kappa^4 \epsilon^{-1})$	$\mathcal{O}(\kappa^3 \epsilon^{-1})$	$\mathcal{O}(\kappa^3 \epsilon^{-1})$
$\text{Gc}(g, \epsilon)$	$\mathcal{O}(\kappa^5 \epsilon^{-5/4})$	$\mathcal{O}(\kappa^4 \epsilon^{-1})$	$\tilde{\mathcal{O}}(\kappa^4 \epsilon^{-1})$
$\text{JV}(g, \epsilon)$	$\mathcal{O}(\kappa^4 \epsilon^{-1})$	$\mathcal{O}(\kappa^3 \epsilon^{-1})$	$\tilde{\mathcal{O}}(\kappa^4 \epsilon^{-1})$
$\text{HV}(g, \epsilon)$	$\mathcal{O}(\kappa^{4.5} \epsilon^{-1})$	$\mathcal{O}(\kappa^{3.5} \epsilon^{-1})$	$\tilde{\mathcal{O}}(\kappa^4 \epsilon^{-1})$

* Require explicit assumption on the sequence.

Table 1: Summary of the convergence results for different algorithms for **deterministic** Euclidean bilevel optimization, including BA (Ghadimi and Wang, 2018), AID-BiO (Ji et al., 2021) and ITD-BiO (Ji et al., 2021). We hide additional $\log(1/\epsilon)$ factors in $\tilde{\mathcal{O}}$.

1.1 Main results

In this work, we first analyze the method of calculating and estimating the hypergradient for bilevel problems on Riemannian manifolds. Our propositions include the Euclidean bilevel

Algorithm	BSA	Stoc-BiO	TTSA*	ALSET	STABLE*
batch size	$\mathcal{O}(1)$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$
y -update	$\mathcal{O}(\epsilon^{-1})$ steps SGD	SGD	SGD	SGD	correction
Gc(F, ϵ)	$\mathcal{O}(\kappa^6 \epsilon^{-2})$	$\mathcal{O}(\kappa^5 \epsilon^{-2})$	$\mathcal{O}(\text{poly}(\kappa) \epsilon^{-2.5})$	$\mathcal{O}(\kappa^5 \epsilon^{-2})$	$\mathcal{O}(\text{poly}(\kappa) \epsilon^{-2})$
Gc(G, ϵ)	$\mathcal{O}(\kappa^9 \epsilon^{-3})$	$\mathcal{O}(\kappa^9 \epsilon^{-2})$	$\mathcal{O}(\text{poly}(\kappa) \epsilon^{-2.5})$	$\mathcal{O}(\kappa^9 \epsilon^{-2})$	$\mathcal{O}(\text{poly}(\kappa) \epsilon^{-2})$
JV(G, ϵ)	$\mathcal{O}(\kappa^6 \epsilon^{-2})$	$\mathcal{O}(\kappa^5 \epsilon^{-2})$	$\mathcal{O}(\text{poly}(\kappa) \epsilon^{-2.5})$	$\mathcal{O}(\kappa^5 \epsilon^{-2})$	$\mathcal{O}(\text{poly}(\kappa) \epsilon^{-2})$
HV(G, ϵ)	$\tilde{\mathcal{O}}(\kappa^6 \epsilon^{-2})$	$\tilde{\mathcal{O}}(\kappa^6 \epsilon^{-2})$	$\tilde{\mathcal{O}}(\text{poly}(\kappa) \epsilon^{-2.5})$	$\tilde{\mathcal{O}}(\kappa^6 \epsilon^{-2})$	$\mathcal{O}(\text{poly}(\kappa) \epsilon^{-2})$

* For algorithms that did not specify the dependence on condition number κ , we use the notation $\text{poly}(\kappa)$ to summarize the κ dependence.

Table 2: Summary of the convergence results for different algorithms for **stochastic** Euclidean bilevel optimization, including BSA (Ghadimi and Wang, 2018), Stoc-BiO (Ji et al., 2021), TTSA (Hong et al., 2023), ALSET (Chen et al., 2021) and STABLE (Chen et al., 2022). For the batch size we only include the ϵ dependency. We hide additional $\log(1/\epsilon)$ factors in $\tilde{\mathcal{O}}$.

problems as special cases and involve the calculation of Riemannian cross derivatives, which are of independent interests to the Riemannian optimization field.

Our contribution also lies in proposing two algorithms (RieBO and RieSBO) for both the problems (1.1) and (1.2) correspondingly. For the deterministic problem (1.1), our analysis shows that with a multi-step inner loop and a single-step outer loop, one could yield the similar gradient complexities Gc(f, ϵ), Gc(g, ϵ), Jacobian- and Hessian-vector product complexities JV(g, ϵ) and HV(g, ϵ) same as the Euclidean counterparts in Ji et al. (2021), as presented in Table 3, as well as for the stochastic problem (1.2) (Chen et al., 2021). It is worth noticing that for the stochastic problem, we adopt the framework of Chen et al. (2021) onto Riemannian manifolds so that the batch size of the hypergradient estimate can be $\mathcal{O}(1)$, significantly smaller than $\mathcal{O}(\epsilon^{-1})$ as in Ji et al. (2021).

Algorithm	RieBO (Algorithm 1)	RieSBO (Algorithm 2)
batch size	No batch	$\mathcal{O}(1)$
y -update	GD	SGD
Gc(F, ϵ)	$\mathcal{O}(\kappa^3 \epsilon^{-1})$	$\mathcal{O}(\kappa^5 \epsilon^{-2})$
Gc(G, ϵ)	$\mathcal{O}(\kappa^4 \epsilon^{-1})$	$\mathcal{O}(\kappa^9 \epsilon^{-2})$
JV(G, ϵ)	$\mathcal{O}(\kappa^3 \epsilon^{-1})$	$\mathcal{O}(\kappa^5 \epsilon^{-2})$
HV(G, ϵ)	$\tilde{\mathcal{O}}(\kappa^{3.5} \epsilon^{-1})$	$\tilde{\mathcal{O}}(\kappa^6 \epsilon^{-2})$

Table 3: Summary of the convergence results for the proposed algorithms in this paper, where all the oracles are with respect to Riemannian gradients and Riemannian second-order derivatives. RieBO (Algorithm 1) solves (1.1), and RieSBO (Algorithm 2) solves (1.2). We hide additional $\log(1/\epsilon)$ factors in $\tilde{\mathcal{O}}$.

Finally, we implement the proposed method in the manifold-constrained bilevel optimization problems, namely the distributionally robust optimization on Riemannian manifolds

with two specific examples: robust maximum likelihood estimation and robust Karcher mean problem on the manifold of positive definite matrices. These numerical results demonstrate the efficiency and potential applicability of the proposed methods.

1.2 Related works

Bilevel optimization. Bilevel optimization problem, also known as nested optimization problem, whose origin dates back to the 50s and 70s (Stackelberg and Peacock, 1952; Bracken and McGill, 1973). Since then, extensive studies have been conducted for solving the bilevel optimization problem (Shi et al., 2005; Moore, 2010). Recently, gradient-based algorithms for solving bilevel optimization problems draw attention because of their applications in machine learning and operations research, such as hyperparameter optimization (Domke, 2012; Pedregosa, 2016; Maclaurin et al., 2015; Franceschi et al., 2018; Lorraine et al., 2020), meta learning (Franceschi et al., 2018; Ji et al., 2021), etc. The increasing attention toward bilevel modeling result in a rich content on the discussion of optimization these bilevel problems, see, e.g. Gould et al. (2016); Shaban et al. (2019); Liu et al. (2020); Li et al. (2020b); Grazzi et al. (2020); Ji and Liang (2023) for some discussion on how to compute the hypergradient (gradient of Φ) and the computational lower bounds. There has also been discussions about the rate of convergence for specific algorithms (Ghadimi and Wang, 2018; Hong et al., 2023; Ji et al., 2021; Chen et al., 2022, 2021). These well-established convergence rate results are summarized in Tables 1 and 2, for deterministic and stochastic problems, respectively. Recently, a line of work (Khanduri et al., 2021; Yang et al., 2023) which utilizes the momentum-based stochastic algorithms can achieve a better oracle complexity of $\mathcal{O}(\epsilon^{-1.5})$ for the Euclidean version of the stochastic problem (1.2). We did not include this line of work since the Riemannian counterparts of these works would rely on the utilization of parallel/vector transport in the algorithm updates. We deliberately avoid these complicated operations in algorithm design and postpone them for future works.

It is worth mentioning that minimax saddle point problems $\min_x \max_y f(x, y)$ are special cases of bilevel optimization problems by taking $g = -f$. Minimax problems are of great interests to the machine learning community (Daskalakis and Panageas, 2018; Mokhtari et al., 2020; Yoon and Ryu, 2021; Lin et al., 2020b). The analysis of this paper relates to the nonconvex-strongly-concave minimax problem on Riemannian manifolds as in Huang and Gao (2023), which showed that the Riemannian gradient descent ascent (RGDA) achieves oracle calls with orders $\mathcal{O}(\kappa^2\epsilon^{-1})$ for the deterministic case and $\mathcal{O}(\kappa^3\epsilon^{-2})$ for the stochastic case. These results match our convergence results in terms of the order of ϵ , but has better κ dependence. This makes sense because our proposed method is a multi- y step GDmax algorithm (see Nouiehed et al. (2019); Jin et al. (2020)) when applied to the minimax problem and naturally has a larger κ dependence. Recently, the authors of Cai et al. (2023) considered the minimax game on Riemannian manifolds under the assumption of geodesic-strongly-monotone (a generalization of strongly-convex-strongly-concave minimax game) and provided a stochastic Riemannian gradient descent-ascent approach which enjoys linear rate of convergence – similar to its Euclidean counterpart. We point out that our work considers nonconvex upper level problems, which is different from the setting in Cai et al. (2023).

Other bilevel-related ongoing research topics include decentralized bilevel optimization (Chen et al., 2024, 2023b; Dong et al., 2023), federate bilevel optimization (Tarzanagh et al., 2022), bilevel without lower strongly convexity (Chen et al., 2023a), to name a few.

Optimization on Riemannian manifolds. Optimization on Riemannian manifolds draws lots of attention recently due to its applications in various fields, including low-rank matrix completion (Boumal and Absil, 2011; Vandereycken, 2013), phase retrieval (Bendory et al., 2017; Sun et al., 2018), dictionary learning (Cherian and Sra, 2016; Sun et al., 2016), dimensionality reduction (Harandi et al., 2017; Tripuraneni et al., 2018; Mishra et al., 2019) and manifold regression (Lin et al., 2017, 2020a). The manifold optimization usually transforms a manifold constrained problem into an unconstrained problem by viewing the manifold as the ambient space and using proper retraction to deal with the loss of linearity, thus achieves better convergence results. For smooth Riemannian optimization, it can be shown that Riemannian gradient descent method require $\mathcal{O}(1/\epsilon)$ iterations to converge to an ϵ -stationary point (i.e. bounding the norm square of the gradient by ϵ) (Boumal et al., 2018). Stochastic algorithms were also studied for smooth Riemannian optimization (Bonnabel, 2013; Zhou et al., 2019; Weber and Sra, 2022; Zhang et al., 2016; Kasai et al., 2018).

The combination of bilevel optimization with Riemannian optimization was largely blank prior to this work. Bonnel et al. (2015) considered a semi-vectorial bilevel optimization model over Riemannian manifolds, which deals with the situation where the lower level problem does not have unique solutions and necessary optimality conditions are provided for their surrogate model. Recently, a concurrent work (Han et al., 2024) also inspect problems in the form of (1.1) and (1.2). We have a very similar algorithm framework with Han et al. (2024) but we would like to point out that, different from Han et al. (2024), our analysis on the stochastic Riemannian bilevel problem (see Theorem 20) is batch-free whereas their result requires a batch size dependent on ϵ (see Han et al. (2024, Theorem 2)). Moreover, we primarily inspect the stochastic Neumann series approximation (4.13) for the stochastic problem (1.2), whereas Han et al. (2024) primarily inspects the approximation methods for the deterministic problem (1.1).

2. Motivating examples

In this section, we provide several motivating examples where at least one between the lower level and the upper level problems are manifold-constrained.

The first set of examples is the robust optimization on Riemannian manifolds, which writes:

$$\min_{y \in \mathcal{N}} \max_{p \in \Delta_n} \sum_{i=1}^n p_i \ell(y; \xi_i) - \lambda \left\| p - \frac{\mathbf{1}}{n} \right\|^2, \quad (2.1)$$

where $\Delta_n := \{y \in \mathbb{R}^n : \sum_{i=1}^n y_i = 1, y_i \geq 0\}$ is the probability simplex, and ℓ is geodesically convex. This problem minimizes n loss functions by dynamically assigning different weights to them, and making sure that the larger loss has larger weights (see Chen et al. (2017); Huang and Gao (2023)). By minimax theorem we can exchange the min and max of the

problem, thus it can be equivalently formulated as a bilevel optimization as follows:

$$\begin{aligned} \min_{p \in \Delta_n} \quad & \lambda \left\| p - \frac{\mathbf{1}}{n} \right\|^2 - \sum_{i=1}^n p_i \ell(y; \xi_i) \\ \text{s.t.} \quad & y \in \operatorname{argmin}_{y \in \mathcal{N}} \sum_{i=1}^n p_i \ell(y; \xi_i). \end{aligned} \tag{2.2}$$

In the numerical experiment section, we inspect two specific examples of robust optimization on Riemannian manifolds.

We also include another example on Riemannian meta-learning (Han et al., 2024). Consider a meta learning problem, where one has m tasks and each task i is represented by a support and query set \mathcal{D}_s^i and \mathcal{D}_q^i , and the target is to learn a set of parameters w which could quickly adapt to all the m tasks while each of the task also has its own training parameter w_i . In Riemannian meta learning, one would require the parameter w lying on the manifold \mathcal{M} , usually the Stiefel manifold (orthogonal constraints). This result in the following Riemannian meta learning problem:

$$\begin{aligned} \min_{w \in \mathcal{M}} \quad & \frac{1}{m} \sum_{i=1}^m \mathcal{L}(w, w_i^*(w); \mathcal{D}_q^i) \\ \text{s.t.} \quad & w_i^*(w) = \operatorname{argmin}_{w_i} \frac{1}{m} \sum_{i=1}^m \mathcal{L}(w, w_i; \mathcal{D}_s^i) + \mathcal{R}(w_i) \end{aligned} \tag{2.3}$$

where \mathcal{L} is a certain loss function that we want to minimize and \mathcal{R} is a regularizer to maintain the strong convexity of the lower-level problem.

We refer to the concurrent work Han et al. (2024) for more examples.

3. Preliminaries on Riemannian Optimization

In this part, we briefly review the basic tools we use for optimization on Riemannian manifolds (Lee, 2006a; Tu, 2011; Boumal, 2023). Suppose \mathcal{M} is an m -dimensional differentiable manifold. The tangent space $T_x \mathcal{M}$ at $x \in \mathcal{M}$ is a linear subspace that consists of the derivatives of all differentiable curves on \mathcal{M} passing through x : $T_x \mathcal{M} := \{\gamma'(0) : \gamma(0) = x, \gamma([- \delta, \delta]) \subset \mathcal{M} \text{ for some } \delta > 0, \gamma \text{ is differentiable}\}$. Notice that for every vector $\gamma'(0) \in T_x \mathcal{M}$, it can be defined in a coordinate-free sense via the operation over smooth functions: $\forall f \in C^\infty(\mathcal{M})$ ¹, $\gamma'(0)(f) := \frac{df \circ \gamma(t)}{dt} |_{t=0}$. The notion of Riemannian manifold is defined as follows.

Definition 1 (Riemannian manifold) *A manifold \mathcal{M} is a Riemannian manifold if it is equipped with an **inner product** on the tangent space. In particular, at any point $x \in \mathcal{M}$,*

1. see Tu (2011, Definition 6.1). For brevity, we omit most of the basic definitions related to differential manifolds and refer to Tu (2011) for details. A smooth mapping on a differential manifold is actually not very different from the Euclidean smooth function since the manifold is locally diffeomorphic to Euclidean spaces.

we have an inner product $\langle \cdot, \cdot \rangle_x : \mathbb{T}_x \mathcal{M} \times \mathbb{T}_x \mathcal{M} \rightarrow \mathbb{R}$, that varies smoothly² on \mathcal{M} . The inner product $\langle \cdot, \cdot \rangle_x$ is usually referred to as the Riemannian metric.

Throughout the paper, we assume that \mathcal{M} and \mathcal{N} are complete Riemannian manifolds equipped with their corresponding Riemannian metrics $\langle \cdot, \cdot \rangle_x$ and $\langle \cdot, \cdot \rangle_y$. Here the completeness refers to the fact that the Riemannian manifold is a complete metric space such that every Cauchy sequence converges. Note that all the matrix manifolds mentioned in this paper are complete.

Note that for the case when $\mathcal{M} = \mathbb{R}^d$ the Euclidean space, the tangent space is $\mathbb{T} \mathcal{M} = \mathbb{T} \mathbb{R}^d$ which is isometric to \mathbb{R}^d , and the Riemannian metric is just the common inner product. Another commonly encountered example is the Stiefel manifold (note that $p = 1$ gives the unit sphere) given by

$$\mathcal{M} = \text{St}(n, p) := \{X \in \mathbb{R}^{n \times p} : X^\top X = I_p\}. \quad (3.1)$$

The tangent space of $\text{St}(n, p)$ is given by $\mathbb{T}_X \mathcal{M} = \{\xi \in \mathbb{R}^{n \times p} : X^\top \xi + \xi^\top X = 0\}$. One could equip the tangent space with common inner product $\langle X, Y \rangle := \text{tr}(X^\top Y)$ as the metric to form a Riemannian manifold. For additional examples, see Absil et al. (2008, Chapter 3) or Boumal (2023, Chapter 7).

Below, we also review the notion of differential of a mapping on manifolds.

Definition 2 (Differential and Riemannian gradients) *Let $F : \mathcal{M} \rightarrow \mathcal{N}$ be a C^∞ map (see Tu (2011, Definition 6.5)) between two differential manifolds. At each point $x \in \mathcal{M}$, the differential of F is a linear mapping (also known as the push-forward):*

$$DF : \mathbb{T}_x \mathcal{M} \rightarrow \mathbb{T}_{F(x)} \mathcal{N},$$

such that $\forall \xi \in \mathbb{T}_x \mathcal{M}$, $DF(\xi) \in \mathbb{T}_x \mathcal{N}$ is given by

$$(DF(\xi))(f) := \xi(f \circ F) \in \mathbb{R}, \quad \forall f \in C^\infty(\mathcal{M}).$$

If $\mathcal{N} = \mathbb{R}$, i.e. $f \in C^\infty(\mathcal{M})$, the differential of f is usually denoted as df . For a Riemannian manifold with Riemannian metric $\langle \cdot, \cdot \rangle$, the Riemannian gradient for $f \in C^\infty(\mathcal{M})$ is the unique tangent vector $\text{grad} f(x) \in \mathbb{T}_x \mathcal{M}$ satisfying

$$df(\xi) = \langle \text{grad} f, \xi \rangle_x, \quad \forall \xi \in \mathbb{T}_x \mathcal{M}.$$

The Riemannian gradient is the generalization of the common Euclidean gradients to Riemannian manifolds, and it is the central concept for our algorithm design. If the manifold is an embedded submanifold of some Euclidean space, then the Riemannian gradient is simply the projection of the Euclidean gradient onto the tangent space. For example, if the manifold is unit sphere $\mathcal{M} = \mathbb{S}^{d-1} = \{x \in \mathbb{R}^d | x^\top x = 1\} \subset \mathbb{R}^d$, then the tangent space at $x \in \mathcal{M}$ is given by $\mathbb{T}_x \mathcal{M} = \{\xi \in \mathbb{R}^d | \xi^\top x = 0\}$, and for any smooth function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, the Riemannian gradient is given by:

$$\text{grad} f(x) = (I_d - xx^\top) \nabla f(x)$$

2. Mapping $\langle \cdot, \cdot \rangle_x$ is actually a function with three inputs on $\mathbb{T} \mathcal{M} \times \mathbb{T} \mathcal{M} \times \mathcal{M}$, and here ‘‘varies smoothly’’ means that it is a smooth mapping on the product manifold $\mathbb{T} \mathcal{M} \times \mathbb{T} \mathcal{M} \times \mathcal{M}$. See Boumal (2023, Definition 5.2) for details.

which is simply projecting $\nabla f(x)$ to the tangent space $T_x \mathcal{M} = \{\xi \in \mathbb{R}^d | \xi^\top x = 0\}$.

For the convergence analysis, we also need the notions of exponential mapping and parallel transport. To this end, we need to first recall the definition of a geodesic.

Definition 3 (Geodesic and exponential mapping) *Given $x \in \mathcal{M}$ and $\xi \in T_x \mathcal{M}$, the geodesic is the curve $\gamma : I \rightarrow \mathcal{M}$, where I is an open set of \mathbb{R} containing 0, such that $\gamma(0) = x$, $\dot{\gamma}(0) = \xi$ and $\nabla_{\dot{\gamma}} \dot{\gamma} = 0$ where $\nabla : T_x \mathcal{M} \times T_x \mathcal{M} \rightarrow T_x \mathcal{M}$ is the Levi-Civita connection defined by metric g . In local coordinate sense, γ is the unique solution of the following second-order differential equations:*

$$\frac{d^2 \gamma^k}{dt^2} + \Gamma_{i,j}^k \frac{d\gamma^i}{dt} \frac{d\gamma^j}{dt} = 0, \quad (3.2)$$

under Einstein summation convention, where $\Gamma_{i,j}^k$ are Christoffel symbols, again defined by metric tensor. The exponential mapping Exp_x is defined as a mapping from $T_x \mathcal{M}$ to \mathcal{M} s.t. $\text{Exp}_x(\xi) := \gamma(1)$ with γ being the geodesic with $\gamma(0) = x$, $\dot{\gamma}(0) = \xi$. A natural corollary is $\text{Exp}_x(t\xi) := \gamma(t)$ for $t \in [0, 1]$. Another useful fact is $\text{dist}(x, \text{Exp}_x(\xi)) = \|\xi\|_x$ since $\gamma'(0) = \xi$ which preserves the speed. Here dist is the geodesic distance which connects the two points by the minimum geodesic.

The definition of the geodesic can be interpreted simply as a notion of “minimum distance curve”. For example, if we solve (3.2) on the unit sphere with Euclidean inner product as the metric, we will simply get curves of great circles starting from one given point and connecting all other points in the shortest distance. For the case of Euclidean space, the Christoffel symbols $\Gamma_{i,j}^k$ will all vanish to zero, and the geodesics are just straight lines (since (3.2) becomes $\gamma''(t) = 0$ and solution is the straight line $\gamma(t) = x_0 + t\xi$).

We want to point out that the definition of the geodesic and exponential mapping is not the main concern of this work, as long as we assume a complete manifold and a minimum geodesic exists between any two points on the manifold. To this end, we always assume that \mathcal{M} is complete throughout this paper, so that Exp_x is always defined for every $\xi \in T_x \mathcal{M}$ (see Lee (2006b, Chapter 6)). For any $x, y \in \mathcal{M}$, the inverse of the exponential mapping $\text{Exp}_x^{-1}(y) \in T_x \mathcal{M}$ is called the logarithm mapping, and we have $\text{dist}(x, y) = \|\text{Exp}_x^{-1}(y)\|_x$, which derives directly from $\text{dist}(x, \text{Exp}_x(\xi)) = \|\xi\|_x$.

Now we give an example of the geodesic and exponential mapping. Consider again the Stiefel manifold (3.1). The geodesic on the Stiefel manifold is given by: $X(t) = [X(0) \quad \dot{X}(0)] \exp\left(t \begin{bmatrix} A(0) & -S(0) \\ I & A(0) \end{bmatrix}\right) \begin{bmatrix} I \\ 0 \end{bmatrix} \exp(-A(0)t)$, for $A(t) = X^\top(t)\dot{X}(t)$ and $S(t) = \dot{X}^\top(t)\dot{X}(t)$ with initial point $X(0)$ and initial speed $\dot{X}(0)$. The exponential mapping is thus given by $\text{Exp}_{X(0)}(\dot{X}(0)) = X(1)$.

With the notion of geodesic, we have the following definition of geodesic convexity and strong convexity, which are the generalizations of their Euclidean counterparts.

Definition 4 (Geodesic (strong) convexity) *A geodesic convex set $\Omega \subset \mathcal{M}$ is a set such that for any two points in the set, there exists a geodesic connecting them that lies entirely in Ω . A function $h : \Omega \rightarrow \mathbb{R}$ is called geodesically convex if for any $p, q \in \Omega$, we have $h(\gamma(t)) \leq (1-t)h(p) + th(q)$, where γ is a geodesic in Ω with $\gamma(0) = p$ and $\gamma(1) = q$. It is called μ -geodesically strongly convex if we have $h(\gamma(t)) \leq (1-t)h(p) + th(q) - \frac{\mu t(1-t)}{2} \text{dist}(p, q)^2$.*

If h is a continuously differentiable function, then it is geodesically convex if and only if (see Boumal (2023, Chapter 11)) $h(y) \geq h(x) + \langle \mathbf{grad}h(x), \mathbf{Exp}_x^{-1}(y) \rangle_x$, and is geodesically strongly convex if and only if $h(y) \geq h(x) + \langle \mathbf{grad}h(x), \mathbf{Exp}_x^{-1}(y) \rangle_x + \frac{\mu}{2} \text{dist}(x, y)^2$.

If h is a twice continuously differentiable function, then it is geodesically convex if and only if (see Boumal (2023, Chapter 11)) $\frac{d^2h(\gamma(t))}{dt^2} \geq 0$, and is geodesically strongly convex if and only if $\frac{d^2h(\gamma(t))}{dt^2} \geq \mu$, where γ is a geodesic.

We also present the definition of parallel transport, which is used in the assumption and the convergence analysis, but not explicitly used in the algorithm updates.

Definition 5 (Parallel transport) Given a Riemannian manifold (\mathcal{M}, g) and two points $x, y \in \mathcal{M}$, the parallel transport $P_{x \rightarrow y} : \mathbb{T}_x \mathcal{M} \rightarrow \mathbb{T}_y \mathcal{M}$ is a linear operator which keeps the inner product: $\forall \xi, \zeta \in \mathbb{T}_x \mathcal{M}$, we have $\langle P_{x \rightarrow y} \xi, P_{x \rightarrow y} \zeta \rangle_y = \langle \xi, \zeta \rangle_x$.

Notice that the existence of parallel transport depends on the curve connecting x and y , which is not a problem for complete Riemannian manifold since we always take the unique geodesic that connects x and y . Parallel transport is useful in our convergence proofs since the Lipschitz condition for the Riemannian gradient requires moving the gradients in different tangent spaces “parallel” to the same tangent space.

For the Euclidean spaces, parallel transport is simply the identity mapping, since the tangent space remains the same at every point. Another example is the Stiefel manifold (3.1), where there is no closed-form expression for the parallel transport, and people usually utilize the projection onto the tangent space, given by $\text{proj}_{\mathbb{T}_X \mathcal{M}}(\xi) = (I - XX^\top)\xi + X \text{skew}(X^\top \xi)$, where $\text{skew}(A) := (A - A^\top)/2$, to transport $\xi \in \mathbb{T}_{X_0} \text{St}(n, p)$ to $\mathbb{T}_X \text{St}(n, p)$. We refer to Absil et al. (2008, Chapter 8) and Boumal (2023, Chapter 10) for additional examples and more discussions on vector and parallel transports.

We also have the following definition of Lipschitz smoothness on the manifolds.

Definition 6 (Geodesic Lipschitz smoothness) A function $h : \Omega \rightarrow \mathbb{R}$ is called geodesic-Lipschitz smooth if we have:

$$\|\mathbf{grad}h(y) - P_{x \rightarrow y} \mathbf{grad}h(x)\| \leq \ell_h \text{dist}(x, y). \quad (3.3)$$

Moreover, we have (see Zhang et al. (2016))

$$h(y) \leq h(x) + \langle \mathbf{grad}h(x), \mathbf{Exp}_x^{-1}(y) \rangle_x + \frac{\ell_h}{2} \text{dist}(x, y)^2. \quad (3.4)$$

Geodesic Lipschitz smoothness is a generalization of the standard gradient-Lipschitz assumption in Euclidean optimization (Nesterov et al., 2018) to the Riemannian setting, and is made in several works (Boumal, 2023; Boumal et al., 2018). To generalize the Euclidean notion to the Riemannian setting, due to the fact that $\mathbf{grad}f(x)$ and $\mathbf{grad}f(y)$ are not in the same tangent space, we need to utilize parallel transports $P_{x \rightarrow y}$ to match the two vectors in the same tangent space.

To proceed to the bilevel hypergradient estimation, we need the notions of Riemannian Hessian and Riemannian cross-derivatives (Jacobians) (see Han et al. (2023)).

Definition 7 (Riemannian Hessian) For function $f : \mathcal{M} \rightarrow \mathbb{R}$, the Riemannian Hessian is a symmetric 2-form $H(f) : \mathbb{T}\mathcal{M} \times \mathbb{T}\mathcal{M} \rightarrow \mathbb{R}$ defined as: $\forall \xi, \eta \in \mathbb{T}\mathcal{M}$,

$$H(f)(\xi, \eta) = \langle \nabla_{\xi} \text{grad} f, \eta \rangle,$$

where ∇ here is the Levi-Civita connection (see Lee (2006a)). H can also be interpreted as a linear map $H(f) : \mathbb{T}\mathcal{M} \rightarrow \mathbb{T}\mathcal{M}$, $\forall \xi \in \mathbb{T}_x \mathcal{M}$,

$$H(f)(\xi) = \nabla_{\xi} \text{grad} f.$$

Note that if f is μ -geodesic strongly convex, we have that $H(f)(\xi, \xi) \geq 0$ with equality attained if and only if $\xi = 0 \in \mathbb{T}\mathcal{M}$ (see Boumal (2023, Chapter 11)).

Definition 8 (Riemannian cross-derivatives) For a smooth function defined on product manifold $f : \mathcal{M} \times \mathcal{N} \rightarrow \mathbb{R}$, the Riemannian cross-derivatives is defined as a linear mapping $\text{grad}_{x,y}^2(f) : \mathbb{T}\mathcal{M} \rightarrow \mathbb{T}\mathcal{N}$ such that $\forall \xi \in \mathbb{T}_x \mathcal{M}$,

$$\text{grad}_{x,y}^2(f)[\xi] = D_x \text{grad}_y f(x, y)[\xi],$$

where D_x is the differential with respect to variable x . $\text{grad}_{y,x}^2(f)$ is defined similarly.

A useful fact is that $\text{grad}_{x,y}^2(f)$ and $\text{grad}_{y,x}^2(f)$ are adjoint operators.

Proposition 9 $\text{grad}_{x,y}^2$ and $\text{grad}_{y,x}^2$ are adjoints, i.e.

$$\langle \eta, \text{grad}_{x,y}^2 f(x, y)[\xi] \rangle_y = \langle \text{grad}_{y,x}^2 f(x, y)[\eta], \xi \rangle_x, \forall \xi \in \mathbb{T}_x \mathcal{M} \text{ and } \forall \eta \in \mathbb{T}_y \mathcal{N},$$

where $f \in \mathcal{C}^1(\mathcal{M})$ is any continuously differentiable function over \mathcal{M} .

Proof Note

$$\langle \eta, \text{grad}_{x,y}^2 f(x, y)[\xi] \rangle_y = \xi(\langle \eta, \text{grad}_y f(x, y) \rangle_y) = \xi(\eta(f)),$$

and similarly

$$\langle \text{grad}_{y,x}^2 f(x, y)[\eta], \xi \rangle_x = \eta(\xi(f)).$$

Note that here ξ and η are actually acting on different coordinates of f . We can extend $\tilde{\xi}(x, y) = (\xi(x), 0) \in \mathbb{T}_x \mathcal{M} \times \mathbb{T}_y \mathcal{N}$ and similarly $\tilde{\eta}(x, y) = (0, \eta(y)) \in \mathbb{T}_x \mathcal{M} \times \mathbb{T}_y \mathcal{N}$. Now subtracting the above equations we have

$$\langle \eta, \text{grad}_{x,y}^2 f(x, y)[\xi] \rangle_y - \langle \text{grad}_{y,x}^2 f(x, y)[\eta], \xi \rangle_x = [\tilde{\xi}, \tilde{\eta}](f),$$

where $[\tilde{\xi}, \tilde{\eta}] = \tilde{\xi}\tilde{\eta} - \tilde{\eta}\tilde{\xi}$ is the Lie bracket. It is easy to verify in local coordinates that $[\tilde{\xi}, \tilde{\eta}]$ is zero since $\tilde{\xi}$ and $\tilde{\eta}$ act on disjoint local coordinates. \blacksquare

4. Bilevel hypergradient estimation on Riemannian manifolds

We first inspect the calculation of the hypergradient for problem (1.1), namely the Riemannian gradient $\text{grad}\Phi(x)$. Notice that y^* is actually a map $\mathcal{M} \rightarrow \mathcal{N}$, thus we need to follow the notion of the differential of maps between manifolds. To calculate the Riemannian gradient $\text{grad}\Phi(x)$, we first need some basic assumptions to ensure the existence and uniqueness of $y^* : \mathcal{M} \rightarrow \mathcal{N}$, also the existence of $\text{grad}\Phi(x)$:

Assumption 1 *The manifolds \mathcal{M} and \mathcal{N} are complete Riemannian manifolds. Moreover, \mathcal{N} is a Hadamard manifold whose sectional curvature is lower bounded by $\iota < 0$ ³.*

We use the notation $\langle \cdot, \cdot \rangle_x$ and $\langle \cdot, \cdot \rangle_y$ to represent their Riemannian metrics, for $x \in \mathcal{M}$ and $y \in \mathcal{N}$. The corresponding norms are $\| \cdot \|_x$ and $\| \cdot \|_y$. Note that from now on we may omit the subscript since the corresponding manifold and tangent space can be identified by the vector in the “.” position.

We also need the following assumptions on the lower and upper objectives:

Assumption 2 (Geodesic strong convexity) *The lower level objective function $g(x, y)$ is second-order continuously differentiable and μ -geodesically strongly convex with respect to y . Note that the total objective $\Phi(x) = f(x, y^*(x))$ may still be nonconvex.*

Such an assumption is necessary in most of the bilevel optimization literature (see for example Ji et al. (2021, Assumption 1)), since we need a unique lower level minimizer $y^*(x)$ for any given $x \in \mathcal{M}$ and the differentiability of $g(x, y)$ will result in a calculable Riemannian gradient $\text{grad}\Phi(x)$. The detail is provided in the proposition right below.

Proposition 10 *Under Assumptions 1 and 2, the mapping $y^* : \mathcal{M} \rightarrow \mathcal{N}$ is differentiable, and the Riemannian gradient $\text{grad}\Phi(x)$ is given by:*

$$\text{grad}\Phi(x) = \text{grad}_x f(x, y^*(x)) - \text{grad}_{y,x}^2 g(x, y^*(x))[v^*(x)], \quad (4.1)$$

where $v^*(x) \in T_{y^*(x)}\mathcal{N}$ is the solution of the following equation:

$$H_y(g(x, y^*(x)))(v) = \text{grad}_y f(x, y^*(x)), \quad (4.2)$$

where H_y is the Riemannian Hessian for the y variable.

Proof We first show that $y^*(x)$ is differentiable. To do this, we first prove that $y^*(x)$ exists and is uniquely determined.

Based on the geodesic strong convexity of g with respect to y , we can show the existence of the solution $y^*(x)$ of the lower level problem by showing that the level sets of function g with respect to y are bounded. Consider $S(x) = \{y \in \mathcal{N} | g(x, y) \leq M\} \neq \emptyset$ and assume that this set is not bounded, then we can find an unbounded sequence $\{y_k\} \subset S(x)$ and without loss of generality we can assume $\|y_k - y_0\| \geq 1, \forall k \geq 1$. Denote $\alpha_k = 1/d(y_k, y_0)$ so that

3. We make this assumption so that the lower function g can be geodesically strongly convex, see Zhang and Sra (2016).

$\lim_k \alpha_k = 0^4$, also denote $z_k = \gamma_k(\alpha_k)$ where γ_k is the geodesic connecting y_0 and y_k with $\gamma_k(0) = y_0$ and $\gamma_k(1) = y_k$. Now based on geodesic strong convexity of g , we have

$$\begin{aligned} g(x, z_k) &\leq \alpha_k g(x, y_k) + (1 - \alpha_k)g(x, y_0) - \frac{\mu(1 - \alpha_k)\alpha_k}{2} \text{dist}(y_k, y_0)^2 \\ &\leq M - \frac{\mu(1 - \alpha_k)\alpha_k}{2} \text{dist}(y_k, y_0)^2 \rightarrow -\infty \end{aligned}$$

which contradicts to the continuity of g . This finishes the existence of the minimizer $y^*(x)$.

Now we proceed to show that the minimizer $y^*(x)$ is unique. For any $y \in \mathcal{N}$ and $\alpha \in (0, 1)$, applying the geodesic strong convexity of g we get:

$$\begin{aligned} \alpha g(x, y) + (1 - \alpha)g(x, y^*(x)) &\geq g(x, \alpha x + (1 - \alpha)y^*(x)) + \frac{\mu\alpha(1 - \alpha)}{2} \text{dist}(y, y^*(x))^2 \\ &\geq g(x, y^*(x)) + \frac{\mu\alpha(1 - \alpha)}{2} \text{dist}(y, y^*(x))^2 \end{aligned}$$

i.e.

$$g(x, y) \geq g(x, y^*(x)) + \frac{\mu(1 - \alpha)}{2} \text{dist}(y, y^*(x))^2$$

which proves the uniqueness of the minimizer $y^*(x)$. We now conclude that $y^* : \mathcal{M} \rightarrow \mathcal{N}$ is a well-defined function.

Next we prove that the mapping $y^* : \mathcal{M} \rightarrow \mathcal{N}$ is differentiable. This is actually directly due to implicit function theorem (see for example, Theorem C.40 in Lee (2006b)). In particular, from the optimality of the lower level problem we know that

$$\mathbf{grad}_y g(x, y) = 0$$

and by strongly convexity of g we know the Hessian $H_y(g)$ is non-singular. Therefore by implicit function theorem, map y^* will be continuously differentiable since we assume g being continuous differentiable.

Now we proceed to calculate the Riemannian gradient $\mathbf{grad}\Phi(x)$ directly. By chain rule,

$$d\Phi(x) = d_x f(x, y^*(x)) + d_y f(x, y^*(x)) \circ (Dy^*(x)), \quad (4.3)$$

where d and D represent the differential operators for real-valued and vector-valued functions, respectively. Notice that the above equation holds in the cotangent space. Since the Riemannian gradients are defined as $\mathbf{grad}\Phi(x) \in T_x \mathcal{M}$ s.t. $\forall \xi \in T_x \mathcal{M}$, $d\Phi(x)(\xi) = \langle \mathbf{grad}\Phi(x), \xi \rangle$, we get from the above equality that

$$\langle \mathbf{grad}\Phi(x), \xi \rangle = \langle \mathbf{grad}_x f(x, y^*(x)), \xi \rangle + \langle \mathbf{grad}_y f(x, y^*(x)), Dy^*(x)(\xi) \rangle. \quad (4.4)$$

Now we have the following optimality condition from the y lower-level problem:

$$\mathbf{grad}_y g(x, y^*(x)) = 0.$$

4. If $d(y_k, y_0) \not\rightarrow \infty$, we can always pick up an subsequence $\{k_j\}$ such that $d(y_{k_j}, y_0) \rightarrow \infty$ due to the unboundedness of $\{y_k\}$.

By taking the differential for x on both sides of the above formula we get: $\forall \xi \in \mathbb{T}_x \mathcal{M}$,

$$\mathbf{grad}_{x,y}^2 g(x, y^*(x))(\xi) + H_y(g(x, y^*(x)))(\mathbf{D}y^*(x)(\xi)) = 0. \quad (4.5)$$

Now taking the inner-product of both sides of the above equation with $v^*(x)$, we get

$$\langle v^*(x), \mathbf{grad}_{x,y}^2 g(x, y^*(x))(\xi) \rangle + \langle \mathbf{grad}_y f(x, y^*(x)), \mathbf{D}y^*(x)(\xi) \rangle = 0.$$

Therefore we get the final result by plugging back the above equation to (4.4) and applying Proposition 9. \blacksquare

When both \mathcal{M} and \mathcal{N} are embedded submanifolds (of two different Euclidean spaces \mathbb{R}^M and \mathbb{R}^N), and $\bar{f} : \mathbb{R}^M \times \mathbb{R}^N \rightarrow \mathbb{R}$ which restricts to $f : \mathcal{M} \times \mathcal{N} \rightarrow \mathbb{R}$ naturally. The Riemannian gradients of f are simply projections of the Euclidean gradients onto the tangent spaces:

$$\mathbf{grad}_x f(x, y) = \mathbf{proj}_{\mathbb{T}_x \mathcal{M}}(\nabla_x \bar{f}(x, y)), \quad \mathbf{grad}_y f(x, y) = \mathbf{proj}_{\mathbb{T}_y \mathcal{N}}(\nabla_y \bar{f}(x, y)), \quad (4.6)$$

and the cross-derivatives are calculated as follows as a matrix:

$$\mathbf{grad}_{x,y}^2 f(x, y) = P_y(\nabla_{x,y}^2 \bar{f}(x, y))P_x, \quad (4.7)$$

where $P_x = \mathbf{proj}_{\mathbb{T}_x \mathcal{M}} \in \mathbb{R}^{M \times M}$ and $P_y = \mathbf{proj}_{\mathbb{T}_y \mathcal{N}} \in \mathbb{R}^{N \times N}$ are projection matrices onto tangent spaces, and $\nabla_{x,y}^2 \bar{f}(x, y) \in \mathbb{R}^{N \times M}$ is the regular partial gradient, namely $[\nabla_{x,y}^2 \bar{f}(x, y)]_{j,i} = \frac{\partial^2 \bar{f}}{\partial y_j \partial x_i}$.

In practice we cannot solve the inner minimization and (4.2) exactly. Suppose we have a point $y \in \mathcal{N}$, we can solve the approximate problem of (4.2):

$$H_y(g(x, y))[v] = \mathbf{grad}_y f(x, y) \quad (4.8)$$

with an N -step conjugate gradient method, yielding $\hat{v}^N(x, y)$, then we can estimate

$$\mathbf{grad} \Phi(x) \approx \mathbf{grad}_x f(x, y) - \mathbf{grad}_{y,x}^2 g(x, y)[\hat{v}^N(x, y)], \quad (4.9)$$

which we further refer to as the approximate implicit differentiation (AID) estimate of (4.1). For the rest of the paper we denote $h_g^{k,t} := \mathbf{grad}_y g(x^k, y^{k,t})$ and

$$h_\Phi(x, y) := \mathbf{grad}_x f(x, y) - \mathbf{grad}_{y,x}^2 g(x, y)[\hat{v}^N(x, y)]. \quad (4.10)$$

We abbreviate the notation by $h_\Phi^k := h_\Phi(x^k, y^k)$ in the algorithms.

For the stochastic problem (1.2), we have an estimate of the Riemannian gradient of the hyperfunction Φ described as follows (see Hong et al. (2023); Chen et al. (2021) for the original Euclidean setting): We first update the inner problem $y^t \leftarrow \mathbf{Exp}_{y^{t-1}}(-\alpha \mathbf{grad}_y G(x, y^{t-1}; \zeta^{t-1}))$ for $t = 1, \dots, T$. Meanwhile the estimate for the gradient of F and the second-order gradient of G will require us to further have independent samples ξ and $\zeta_{(q)}$, $q = 0, 1, \dots, Q$, so that we define the stochastic gradient estimator as:

$$\mathbf{grad} \Phi(x) \approx \mathbf{grad}_x F(x, y^T; \xi) - \mathbf{grad}_{y,x}^2 G(x, y^T; \zeta_0)[v_Q(x, y^T)], \quad (4.11)$$

where v_Q is the approximation of (4.2), defined as (see Hong et al. (2023); Liao et al. (2018)):

$$v_Q(x, y) := \eta Q \prod_{q=1}^{Q'} (I - \eta H_y(G(x, y; \zeta_{(q)}))) [\text{grad}_y F(x, y; \xi)], \quad (4.12)$$

where Q' is drawn uniformly from $\{0, 1, \dots, Q - 1\}$ and the extra parameter η will be later determined to ensure a better approximation to (4.2), motivated by the Neumann series $\sum_{i=0}^{\infty} U^i = (I - U)^{-1}$.

From now on we denote $\tilde{h}_g^{k,t} = \text{grad}_y G(x^k, y^{k,t}; \zeta_{k,t})$ and

$$\tilde{h}_\Phi^k := \text{grad}_x F(x^k, y^k; \xi_k) - \text{grad}_{y,x}^2 G(x^k, y^k; \zeta_{k,(0)}) [v_Q^k], \quad (4.13)$$

where

$$v_Q^k := \eta Q \prod_{q=1}^{Q'} (I - \eta H_y(G(x^k, y^k; \zeta_{k,(q)}))) [\text{grad}_y F(x^k, y^k; \xi_k)].$$

5. Deterministic Algorithm RieBO and Its Convergence

We propose RieBO (Algorithm 1) for the deterministic bilevel manifold optimization (1.1). The algorithm is a generalization of its Euclidean counterpart proposed in Ji et al. (2021) where we employ conjugate gradient method to solve the hypergradient estimation problem (4.10).

For the deterministic case, we utilize the following notion of stationarity:

Definition 11 *A point $x \in \mathcal{M}$ is called an ϵ -stationary point for (1.1) if $\|\text{grad}\Phi(x)\|^2 \leq \epsilon$.*

Following Zhang and Sra (2016), it is very important that the quantity $\tau(\iota, \text{dist}(y^{k,t}, y^*(x^k)))$ is bounded during the lower level update, where

$$\tau(\iota, c) := \frac{\sqrt{|\iota|c}}{\tanh(\sqrt{|\iota|c})}. \quad (5.1)$$

Therefore we need the following assumption along with aforementioned assumptions.

Assumption 3 *We assume that the quantity $\tau(\iota, \text{dist}(y^{k,t}, y^*(x^k)))$ is always upper bounded by τ for all k and t ⁵.*

Assumption 3 is a commonly used assumption for convex optimization on Riemannian manifolds (see e.g. Zhang and Sra (2016, Corollary 8)). Geometrically, τ quantifies how twisted distance between the update sequence and the optimal point becomes when the manifold has a curvature ι . For Euclidean space, this quantity τ is always zero, and the farther the manifold is from Euclidean space, the larger this quantity becomes.

The next assumptions are standard in manifold optimization literature about the Lipschitz smoothness.

5. Note that this assumption is satisfied if the lower level problem is conducted in a compact subset in \mathcal{N} and assuming that the iterates of the algorithm stay in this compact region, see Zhang and Sra (2016).

Assumption 4 (Lipschitz Smoothness) For simplicity denote $z = (x, y)$ and $z' = (x', y')$, also denote $\text{dist}(z, z') = \sqrt{\text{dist}(x, x')^2 + \text{dist}(y, y')^2}$ (note that these two distances are on different manifolds). Moreover, f and g satisfy the following assumptions.

- f satisfies $\ell_{f,0}$ -Lipschitzness:

$$|f(z) - f(z')| \leq \ell_{f,0} \text{dist}(z, z').$$

- $\text{grad}f = [\text{grad}_x f, \text{grad}_y f]$ and $\text{grad}g = [\text{grad}_x g, \text{grad}_y g]$ are $\ell_{f,1}$ and $\ell_{g,1}$ Lipschitz, i.e.,

$$\|\text{grad}f(z) - P_{z' \rightarrow z} \text{grad}f(z')\| \leq \ell_{f,1} \text{dist}(z, z')$$

$$\|\text{grad}g(z) - P_{z' \rightarrow z} \text{grad}g(z')\| \leq \ell_{g,1} \text{dist}(z, z'),$$

where $P_{z' \rightarrow z}$ is the parallel transport on the product manifold $\mathcal{M} \times \mathcal{N}$ and the norm on the left hand side is also induced by the product Riemannian metric on $\mathcal{M} \times \mathcal{N}$ ⁶.

Assumption 4 is a generalization of relatively standard assumptions from bilevel optimization (see e.g. Hong et al. (2023, Assumption 1 and 2)). Note that parallel transports reduces to identity mapping for Euclidean spaces, and the above assumption is exactly the standard Lipschitz continuous and Lipschitz smooth assumptions. Further, for bilevel optimization, most of the works (see e.g. Ji et al. (2021, Assumption 3) and Hong et al. (2023, Assumption 2)) would require the Hessian of the lower level problem also being Lipschitz continuous, in order to make sure the Riemannian hypergradient $\text{grad}\Phi(x)$ Lipschitz smooth. Therefore we need the following assumption.

Assumption 5 (Hessian Smoothness) The second-order derivatives $\text{grad}_{x,y}^2 g(z)$ and $H_y(g(z))$ are $\ell_{g,2}$ Lipschitz (we use the same constant here for simplicity), i.e. for $z = (x, y)$ and $z' = (x', y')$, we have

$$\|\text{grad}_{x,y}^2 g(z) - P_{y^*(x') \rightarrow y^*(x)} \circ \text{grad}_{x,y}^2 g(z') \circ P_{x' \rightarrow x}\|_{\text{op}} \leq \ell_{g,2} \text{dist}(z, z')$$

$$\|H_y(g(z)) - P_{y^*(x') \rightarrow y^*(x)} \circ H_y(g(z')) \circ P_{y^*(x) \rightarrow y^*(x')}\|_{\text{op}} \leq \ell_{g,2} \text{dist}(z, z').$$

Here the norms on the left hand side are the operator norms. We will keep the subscript op whenever it comes to the operator norm, in order to distinguish it from the norms on the tangent spaces.

Now we are ready to conduct our convergence analysis. We first have the following smoothness lemma under the above assumptions.

Lemma 12 Suppose Assumptions 1, 2, 3, 4 and 5 hold, then functions $y^*(x)$ and $\Phi(x) := f(x, y^*(x))$ satisfy: $\forall x, x' \in \mathcal{M}$

$$\text{dist}(y^*(x), y^*(x')) \leq \kappa \text{dist}(x, x'), \quad \kappa = \max \left\{ \frac{\ell_{g,1}}{\mu}, \frac{\ell_{g,2}}{\mu} \right\} \quad (5.2a)$$

6. It is worth noticing that in our convergence result, we need $\tau < \ell_{g,1}/2$. We comment here that this assumption is not a big issue since if g is Lipschitz smooth with parameter $\ell_{g,1}$, then it is also Lipschitz smooth with any parameters $\ell > \ell_{g,1}$. We could always pick up a parameter that satisfies $\tau < \ell_{g,1}/2$, with some sacrifice of the convergence speed. Note that in the Euclidean case $\tau = 0$ so $\tau < \ell_{g,1}/2$ holds naturally.

$$\|Dy^*(x) - P_{y^*(x') \rightarrow y^*(x)} \circ Dy^*(x') \circ P_{x \rightarrow x'}\|_{\text{op}} \leq L_{y^*} \text{dist}(x, x') \quad (5.2b)$$

$$\|\text{grad}\Phi(x) - P_{x' \rightarrow x} \text{grad}\Phi(x')\| \leq L_{\Phi} \text{dist}(x, x'), \quad (5.2c)$$

where

$$L_{y^*} := \left(1 + \frac{\ell_{g,2}}{\mu}\right) \frac{\ell_{g,2}}{\mu} \sqrt{1 + \kappa^2} = \mathcal{O}(\kappa^3), \quad (5.3)$$

and

$$L_{\Phi} := \ell_{f,1} \sqrt{1 + \kappa^2} + \ell_{g,2} \frac{\ell_{f,0}}{\mu} + \ell_{g,1} \left(\frac{\ell_{f,0} \ell_{g,2}}{\mu^2} \sqrt{1 + \kappa^2} + \frac{\ell_{f,1}}{\mu} \right) = \mathcal{O}(\kappa^3). \quad (5.4)$$

Proof For (5.2a), by (4.5) and our Assumptions 3, 4 and 5, we have

$$\|Dy^*(x)\|_{\text{op}} = \|(H_y(g(x, y^*(x))))^{-1} \circ \text{grad}_{x,y}^2 g(x, y^*(x))\|_{\text{op}} \leq \frac{\ell_{g,1}}{\mu}.$$

We thus obtain (5.2a) by a mean value theorem argument in local coordinates (see this link for a detailed proof).

For (5.2b), we have

$$\begin{aligned} & \|Dy^*(x) - P_{y^*(x') \rightarrow y^*(x)} \circ Dy^*(x') \circ P_{x' \rightarrow x}\| \\ &= \|(H_y(g(x, y^*(x))))^{-1} \circ \text{grad}_{x,y}^2 g(x, y^*(x)) \\ &\quad - P_{y^*(x') \rightarrow y^*(x)} \circ (H_y(g(x', y^*(x'))))^{-1} \circ \text{grad}_{x',y'}^2 g(x', y^*(x')) \circ P_{x' \rightarrow x}\| \\ &\leq \|(H_y(g(x, y^*(x))))^{-1} \circ \text{grad}_{x,y}^2 g(x, y^*(x)) \\ &\quad - (H_y(g(x, y^*(x))))^{-1} \circ P_{y^*(x') \rightarrow y^*(x)} \circ \text{grad}_{x,y}^2 g(x', y^*(x')) \circ P_{x' \rightarrow x}\| \\ &\quad + \|(H_y(g(x, y^*(x))))^{-1} \circ P_{y^*(x') \rightarrow y^*(x)} \circ \text{grad}_{x,y}^2 g(x', y^*(x')) \\ &\quad - P_{y^*(x') \rightarrow y^*(x)} \circ (H_y(g(x', y^*(x'))))^{-1} \circ \text{grad}_{x',y'}^2 g(x', y^*(x'))\| \\ &\leq \|(H_y(g(x, y^*(x))))^{-1}\| \|\text{grad}_{x,y}^2 g(x, y^*(x)) - P_{y^*(x') \rightarrow y^*(x)} \circ \text{grad}_{x,y}^2 g(x', y^*(x')) \circ P_{x' \rightarrow x}\| \\ &\quad + \|(H_y(g(x, y^*(x))))^{-1} - P_{y^*(x') \rightarrow y^*(x)} \circ (H_y(g(x', y^*(x'))))^{-1} \circ P_{y^*(x) \rightarrow y^*(x')}\| \|\text{grad}_{x,y}^2 g(x', y^*(x'))\| \\ &\leq \frac{\ell_{g,2}}{\mu} \sqrt{\text{dist}(x, x')^2 + \text{dist}(y^*(x), y^*(x'))^2} \\ &\quad + \ell_{g,1} \|(H_y(g(x, y^*(x))))^{-1} - P_{y^*(x') \rightarrow y^*(x)} \circ (H_y(g(x', y^*(x'))))^{-1} \circ P_{y^*(x) \rightarrow y^*(x')}\| \\ &\leq \frac{\ell_{g,2}}{\mu} \sqrt{1 + \kappa^2} \text{dist}(x, x') \\ &\quad + \ell_{g,1} \|(H_y(g(x, y^*(x))))^{-1} - P_{y^*(x') \rightarrow y^*(x)} \circ (H_y(g(x', y^*(x'))))^{-1} \circ P_{y^*(x) \rightarrow y^*(x')}\|_{\text{op}} \end{aligned}$$

where in the second last inequality we used Assumptions 3, 4 and 5, and we used (5.2a) for the last inequality. Denote $H_1 = H_y(g(x, y^*(x)))$, $P = P_{y^*(x') \rightarrow y^*(x)}$ so that $P_{y^*(x) \rightarrow y^*(x')} = P^{-1}$ and $H_2 = H_y(g(x', y^*(x')))$, then the last term in the above formula becomes:

$$\begin{aligned} & \|H_1^{-1} - P H_2^{-1} P\|_{\text{op}} = \|H_1^{-1} P^{-1} (H_2 - P^{-1} H_1 P) H_2^{-1} P^{-1}\|_{\text{op}} \\ &\leq \frac{1}{\mu^2} \|H_2 - P^{-1} H_1 P\|_{\text{op}} \leq \frac{\ell_{g,2}}{\mu^2} \sqrt{\text{dist}(x, x')^2 + \text{dist}(y^*(x), y^*(x'))^2} \\ &\leq \frac{\ell_{g,2}}{\mu^2} \sqrt{1 + \kappa^2} \text{dist}(x, x'). \end{aligned} \quad (5.5)$$

Here we used Assumptions 3, 5, (5.2a) and the fact that the parallel transport $P = P_{y^*(x') \rightarrow y^*(x)}$ is an isometry, i.e., $\|P\|_{\text{op}} = 1$. Plugging this back we get

$$\|Dy^*(x) - P_{y^*(x') \rightarrow y^*(x)} \circ Dy^*(x') \circ P_{x' \rightarrow x}\|_{\text{op}} \leq \left(\frac{\ell_{g,2}}{\mu} + \frac{\ell_{g,1}\ell_{g,2}}{\mu^2} \right) \sqrt{1 + \kappa^2} \text{dist}(x, x'),$$

which gives (5.2b).

We now show (5.2c). Using (4.1), we get

$$\begin{aligned} \|\text{grad}\Phi(x) - P_{x' \rightarrow x} \text{grad}\Phi(x')\| &\leq \|\text{grad}_x f(x, y^*(x)) - P_{x' \rightarrow x} \text{grad}_x f(x', y^*(x'))\| \\ &+ \|\text{grad}_{y,x}^2 g(x, y^*(x))[v^*(x)] - P_{x' \rightarrow x} \text{grad}_{y,x}^2 g(x', y^*(x'))[v^*(x')]\|. \end{aligned} \quad (5.6)$$

For the first term on the right hand side of (5.6), by Assumption 4 and (5.2a), we have

$$\begin{aligned} &\|\text{grad}_x f(x, y^*(x)) - P_{x' \rightarrow x} \text{grad}_x f(x', y^*(x'))\| \\ &\leq \ell_{f,1} \sqrt{\text{dist}(x, x')^2 + \text{dist}(y^*(x), y^*(x'))^2} \leq \ell_{f,1} \sqrt{1 + \kappa^2} \text{dist}(x, x'). \end{aligned}$$

For the second term on the right hand side of (5.6), we have

$$\begin{aligned} &\|\text{grad}_{y,x}^2 g(x, y^*(x))[v^*(x)] - P_{x' \rightarrow x} \text{grad}_{y,x}^2 g(x', y^*(x'))[v^*(x')]\| \\ = &\|\text{grad}_{y,x}^2 g(x, y^*(x))[v^*(x)] - P_{x' \rightarrow x} \text{grad}_{y,x}^2 g(x', y^*(x')) \left[P_{x \rightarrow x'} P_{x' \rightarrow x} [v^*(x')] \right]\| \\ \leq &\|\text{grad}_{y,x}^2 g(x, y^*(x))[v^*(x)] - P_{x' \rightarrow x} \text{grad}_{y,x}^2 g(x', y^*(x')) P_{x \rightarrow x'} [v^*(x)]\| \\ &+ \|P_{x' \rightarrow x} \text{grad}_{y,x}^2 g(x', y^*(x')) P_{x \rightarrow x'} [v^*(x)] - P_{x' \rightarrow x} \text{grad}_{y,x}^2 g(x', y^*(x')) P_{x \rightarrow x'} P_{x' \rightarrow x} [v^*(x')]\| \\ \leq &\|\text{grad}_{y,x}^2 g(x, y^*(x)) - P_{x' \rightarrow x} \text{grad}_{y,x}^2 g(x', y^*(x')) P_{x \rightarrow x'}\|_{\text{op}} \|v^*(x)\| \\ &+ \|\text{grad}_{y,x}^2 g(x', y^*(x'))\|_{\text{op}} \|v^*(x) - P_{x' \rightarrow x} v^*(x')\|. \end{aligned}$$

Since $v^*(x)$ is the solution of (4.2), also by Assumptions 3 and 4, we have $\|v^*(x)\| \leq \ell_{f,0}/\mu$, also

$$\begin{aligned} &\|v^*(x) - P_{x' \rightarrow x} v^*(x')\| \\ = &\|(H_y(g(x, y^*(x))))^{-1} \text{grad}_y f(x, y^*(x)) - P_{x' \rightarrow x} (H_y(g(x', y^*(x'))))^{-1} \text{grad}_y f(x', y^*(x'))\| \\ = &\|(H_y(g(x, y^*(x))))^{-1} \text{grad}_y f(x, y^*(x)) - P_{x' \rightarrow x} (H_y(g(x', y^*(x'))))^{-1} P_{x \rightarrow x'} P_{x' \rightarrow x} \text{grad}_y f(x', y^*(x'))\| \\ \leq &\|(H_y(g(x, y^*(x))))^{-1} - P_{x' \rightarrow x} (H_y(g(x', y^*(x'))))^{-1} P_{x \rightarrow x'}\|_{\text{op}} \|\text{grad}_y f(x, y^*(x))\| \\ &+ \|(H_y(g(x', y^*(x'))))^{-1}\|_{\text{op}} \|\text{grad}_y f(x, y^*(x)) - P_{x' \rightarrow x} \text{grad}_y f(x', y^*(x'))\| \\ \leq &\left(\frac{\ell_{f,0}\ell_{g,2}}{\mu^2} \sqrt{1 + \kappa^2} + \frac{\ell_{f,1}}{\mu} \right) \text{dist}(x, x'), \end{aligned} \quad (5.7)$$

where we used (5.5), Assumptions 3 and 4.

Combining the above bounds and plugging it to (5.6), we get

$$\begin{aligned} &\|\text{grad}\Phi(x) - P_{x' \rightarrow x} \text{grad}\Phi(x')\| \\ \leq &\ell_{f,1} \sqrt{1 + \kappa^2} \text{dist}(x, x') + \ell_{g,2} \frac{\ell_{f,0}}{\mu} \text{dist}(x, x') + \ell_{g,1} \left(\frac{\ell_{f,0}\ell_{g,2}}{\mu^2} \sqrt{1 + \kappa^2} + \frac{\ell_{f,1}}{\mu} \right) \text{dist}(x, x'), \end{aligned}$$

Algorithm 1: Algorithm for **Riemannian** (deterministic) **Bilevel Optimization** (**RieBO**)

input : K, T, N (steps for conjugate gradient), stepsize $\{\alpha_k, \beta_k\}$, initializations $x^0 \in \mathcal{M}, y^0 \in \mathcal{N}$
for $k = 0, 1, 2, \dots, K - 1$ **do**
 Set $y^{k,0} = y^{k-1}$;
 for $t = 0, \dots, T - 1$ **do**
 Update $y^{k,t+1} \leftarrow \text{Exp}_{y^{k,t}}(-\beta_k h_g^{k,t})$ with $h_g^{k,t} := \text{grad}_y g(x^k, y^{k,t})$;
 end
 Set $y^k \leftarrow y^{k,T}$;
 Update $x^{k+1} \leftarrow \text{Exp}_{x^k}(-\alpha_k h_\Phi^k)$ as in (4.10), where $\hat{v}^N(x^k, y^k)$ is given by an N -step conjugate gradient update, with $\hat{v}^0(x^k, y^k) = P_{y^{k-1} \rightarrow y^k} \hat{v}^N(x^{k-1}, y^{k-1})$;
end

which proves (5.2c). ■

Now we are ready to provide our convergence analysis result for **RieBO** (Algorithm 1). This result is an extension of the convergence result in Ji et al. (2021) to the Riemannian setting.

Theorem 13 *Suppose Assumptions 1, 2, 3, 4 and 5 hold, and take the parameters $\beta_k = \beta \leq \frac{1}{\ell_{g,1}}$, $\alpha_k = \alpha \leq \frac{1}{8L_\Phi}$, $T \geq \mathcal{O}(\kappa)$ and conjugate gradient iteration number $N \geq \mathcal{O}(\sqrt{\kappa})$. Then **RieBO** (Algorithm 1) satisfies:*

$$\frac{1}{K} \sum_{k=0}^{K-1} \|\text{grad}\Phi(x^k)\|^2 \leq \mathcal{O}\left(\frac{L_\Phi}{K}\right). \quad (5.8)$$

The specific choice parameters are given in the proof for the simplicity of the statement. In order to achieve an ϵ -accurate stationary point, the complexity is given by:

- *Gradients:* $\text{Gc}(f, \epsilon) = \mathcal{O}(\kappa^3 \epsilon^{-1})$, $\text{Gc}(g, \epsilon) = \mathcal{O}(\kappa^4 \epsilon^{-1})$;
- *Jacobian and Hessian-vector products:* $\text{JV}(g, \epsilon) = \mathcal{O}(\kappa^3 \epsilon^{-1})$, $\text{HV}(g, \epsilon) = \mathcal{O}(\kappa^{3.5} \epsilon^{-1})$.

To prove this theorem, we need the following lemmas. The first lemma quantifies the error when optimizing (4.8) with N -step conjugate gradient method, see Grazzi et al. (2020, Equation (17))⁷.

7. Note that here the Hessian matrix $H_y(g(x, y))$ is full-rank in the tangent space. However if it is an embedded submanifold then $H_y(g(x, y))$ is actually rank-deficient matrix in the ambient Euclidean space. This is not a concern for showing the linear rate of convergence since we can always conduct CG steps only on the tangent spaces (as Euclidean subspaces of the ambient Euclidean space). It is known that the convergence is still linear even if $H_y(g(x, y))$ is rank-deficient, see Hayami (2018) for a detailed inspection.

Lemma 14 *Suppose we solve (4.8) with N -step conjugate gradient method with the initial point $\hat{v}^0(x, y)$ and output $\hat{v}^N(x, y)$, then we have*

$$\|\hat{v}^N(x, y) - \tilde{v}\| \leq \sqrt{\kappa} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^N \|\hat{v}^0(x, y) - \tilde{v}\|,$$

where \tilde{v} is the exact solution of (4.8).

The next lemma quantifies the error of the inner loop, i.e. the T steps where we do Riemannian gradient descent for the lower problem in RieBO (Algorithm 1).

Lemma 15 *Suppose Assumptions 1, 2, 3, 4 and 5 hold, and we take $\beta_k = \beta = 1/\ell_{g,1}$ as a constant, then RieBO satisfies:*

$$\text{dist}(y^{k,T}, y^*(x^k))^2 \leq (1 - 2\mu\tau\beta^2)^T \text{dist}(y^{k,0}, y^*(x^k))^2. \quad (5.9)$$

Proof For simplicity, we denote $h(y) = g(x^k, y)$, so that $y^*(x^k)$ is the optimal solution of h . We also omit k in this proof, i.e., the update becomes:

$$y^{t+1} \leftarrow \text{Exp}_{y^t}(-\beta_k \text{grad}h(y^t)).$$

By the notions of geodesic smoothness and geodesic convexity, we have

$$\begin{aligned} h(y^{t+1}) - h(y^*) &= h(y^{t+1}) - h(y^t) + h(y^t) - h(y^*) \\ &\leq \langle \text{grad}h(y^t), \text{Exp}_{y^t}(y^{t+1}) \rangle + \frac{\ell_{g,1}}{2} \|\text{Exp}_{y^t}(y^{t+1})\|^2 - \langle \text{grad}h(y^t), \text{Exp}_{y^t}(y^*) \rangle - \frac{\mu}{2} \text{dist}(y^t, y^*)^2 \\ &= -\left(\beta - \frac{\beta^2 \ell_{g,1}}{2}\right) \|\text{grad}h(y^t)\|^2 - \langle \text{grad}h(y^t), \text{Exp}_{y^t}(y^*) \rangle - \frac{\mu}{2} \text{dist}(y^t, y^*)^2, \end{aligned}$$

i.e.,

$$\left(\beta - \frac{\beta^2 \ell_{g,1}}{2}\right) \|\text{grad}h(y^t)\|^2 - \langle \text{grad}h(y^t), \text{Exp}_{y^t}(y^*) \rangle \leq -\frac{\mu}{2} \text{dist}(y^t, y^*)^2. \quad (5.10)$$

Now by Zhang and Sra (2016, Corollary 8), we have,

$$\begin{aligned} \text{dist}(y^{t+1}, y^*)^2 &\leq \text{dist}(y^t, y^*)^2 + 2\beta \langle \text{grad}h(y^t), \text{Exp}_{y^t}(y^*) \rangle + \tau\beta^2 \|\text{grad}h(y^t)\|^2 \\ &\leq (1 - 2\mu\tau\beta^2) \text{dist}(y^t, y^*)^2, \end{aligned}$$

where the last inequality is by (5.10) and $\beta = 1/\ell_{g,1}$. The proof is done by repeatedly applying the above inequality from $t = T - 1$ back to $t = 0$. \blacksquare

The next lemma quantifies the error between our estimation h_{Φ}^k and the true upper level gradient $\text{grad}\Phi(x^k)$.

Lemma 16 *Suppose Assumptions 1, 2, 3, 4 and 5 hold, then RieBO satisfies:*

$$\begin{aligned} \|h_{\Phi}^k - \text{grad}\Phi(x^k)\| &\leq \Gamma(1 - 2\mu\tau\beta^2)^{T/2} \text{dist}(y^*(x^k), y^{k-1}) \\ &\quad + \ell_{g,1} \sqrt{\kappa} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^N \|\hat{v}^0(x^k, y^k) - P_{y^*(x^k) \rightarrow y^k} v^*(x^k)\|, \end{aligned} \quad (5.11)$$

where h_{Φ}^k is the estimate from (4.10) and we have the parameters:

$$\begin{aligned} \tilde{v}^k &= (H_y(g(x^k, y^k)))^{-1} \mathbf{grad}_y f(x^k, y^k) \\ \Gamma &= \ell_{f,1} + \frac{\ell_{f,0} \ell_{g,2}}{\mu} + \ell_{g,1} \left(1 + \sqrt{\kappa} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^N \right) \left(\ell_{g,2} \ell_{f,0} + \frac{\ell_{f,1}}{\mu} \right). \end{aligned} \quad (5.12)$$

Proof We first restate the expression (4.1) for $\mathbf{grad}\Phi(x)$ and (4.10) for h_{Φ}^k :

$$\begin{aligned} \mathbf{grad}\Phi(x^k) &= \mathbf{grad}_x f(x^k, y^*(x^k)) - \mathbf{grad}_{y,x}^2 g(x^k, y^*(x^k)) [v^*(x^k)], \\ h_{\Phi}(x^k, y^k) &:= \mathbf{grad}_x f(x^k, y^k) - \mathbf{grad}_{y,x}^2 g(x^k, y^k) [\hat{v}^N(x^k, y^k)]. \end{aligned}$$

Thus,

$$\begin{aligned} \|h_{\Phi}^k - \mathbf{grad}\Phi(x^k)\| &\leq \|\mathbf{grad}_x f(x^k, y^*(x^k)) - \mathbf{grad}_x f(x^k, y^k)\| \\ &\quad + \|\mathbf{grad}_{y,x}^2 g(x^k, y^*(x^k)) [v^*(x^k)] - \mathbf{grad}_{y,x}^2 g(x^k, y^k) [\hat{v}^N(x^k, y^k)]\| \\ &\leq \|\mathbf{grad}_x f(x^k, y^*(x^k)) - \mathbf{grad}_x f(x^k, y^k)\| \\ &\quad + \|\mathbf{grad}_{y,x}^2 g(x^k, y^*(x^k)) [v^*(x^k)] - \mathbf{grad}_{y,x}^2 g(x^k, y^k) [P_{y^*(x^k) \rightarrow y^k} v^*(x^k)]\| \\ &\quad + \|\mathbf{grad}_{y,x}^2 g(x^k, y^k) [P_{y^*(x^k) \rightarrow y^k} v^*(x^k)] - \mathbf{grad}_{y,x}^2 g(x^k, y^k) [\hat{v}^N(x^k, y^k)]\| \\ &\leq \ell_{f,1} \text{dist}(y^*(x^k), y^k) \\ &\quad + \|\mathbf{grad}_{y,x}^2 g(x^k, y^*(x^k)) - \mathbf{grad}_{y,x}^2 g(x^k, y^k) \circ P_{y^*(x^k) \rightarrow y^k}\|_{\text{op}} \|v^*(x^k)\| \\ &\quad + \|\mathbf{grad}_{y,x}^2 g(x^k, y^k)\|_{\text{op}} \|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \hat{v}^N(x^k, y^k)\| \\ &\leq \left(\ell_{f,1} + \frac{\ell_{f,0} \ell_{g,2}}{\mu} \right) \text{dist}(y^*(x^k), y^k) + \ell_{g,1} \|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \hat{v}^N(x^k, y^k)\|. \end{aligned}$$

Following Lemma 14, we have

$$\begin{aligned} \|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \hat{v}^N(x^k, y^k)\| &\leq \|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \tilde{v}^k\| + \|\tilde{v}^k - \hat{v}^N(x^k, y^k)\| \\ &\leq \|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \tilde{v}^k\| + \sqrt{\kappa} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^N \|\hat{v}^0(x^k, y^k) - \tilde{v}^k\| \\ &\leq \left(1 + \sqrt{\kappa} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^N \right) \|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \tilde{v}^k\| + \sqrt{\kappa} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^N \|\hat{v}^0(x^k, y^k) - P_{y^*(x^k) \rightarrow y^k} v^*(x^k)\|. \end{aligned}$$

For $\|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \tilde{v}^k\|$, by the definitions of \tilde{v}_k and v_k^* , we have

$$\begin{aligned} &\|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \tilde{v}^k\| \\ &= \|P_{y^*(x^k) \rightarrow y^k} (H_y(g(x^k, y^*(x^k))))^{-1} \mathbf{grad}_y f(x^k, y^*(x^k)) - (H_y(g(x^k, y^k)))^{-1} \mathbf{grad}_y f(x^k, y^k)\| \\ &\leq \|P_{y^*(x^k) \rightarrow y^k} (H_y(g(x^k, y^*(x^k))))^{-1} \mathbf{grad}_y f(x^k, y^*(x^k)) - (H_y(g(x^k, y^k)))^{-1} P_{y^*(x^k) \rightarrow y^k} \mathbf{grad}_y f(x^k, y^*(x^k))\| \\ &\quad + \|(H_y(g(x^k, y^k)))^{-1} P_{y^*(x^k) \rightarrow y^k} \mathbf{grad}_y f(x^k, y^*(x^k)) - (H_y(g(x^k, y^k)))^{-1} \mathbf{grad}_y f(x^k, y^k)\| \\ &\leq \|(H_y(g(x^k, y^*(x^k))))^{-1} - P_{y^k \rightarrow y^*(x^k)} (H_y(g(x^k, y^k)))^{-1} P_{y^*(x^k) \rightarrow y^k}\|_{\text{op}} \|\mathbf{grad}_y f(x^k, y^*(x^k))\| \\ &\quad + \|(H_y(g(x^k, y^k)))^{-1}\|_{\text{op}} \|P_{y^*(x^k) \rightarrow y^k} \mathbf{grad}_y f(x^k, y^*(x^k)) - \mathbf{grad}_y f(x^k, y^k)\| \\ &\leq \left(\ell_{g,2} \ell_{f,0} + \frac{\ell_{f,1}}{\mu} \right) \text{dist}(y^k, y^*(x^k)). \end{aligned} \quad (5.13)$$

Therefore, we get

$$\begin{aligned}
 \|h_{\Phi}^k - \text{grad}\Phi(x^k)\| &\leq \left(\ell_{f,1} + \frac{\ell_{f,0}\ell_{g,2}}{\mu} \right) \text{dist}(y^*(x^k), y^k) \\
 &+ \ell_{g,1} \left(1 + \sqrt{\kappa} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^N \right) \|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \tilde{v}^k\| \\
 &+ \ell_{g,1} \sqrt{\kappa} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^N \|\hat{v}^0(x^k, y^k) - P_{y^*(x^k) \rightarrow y^k} v^*(x^k)\| \\
 &\leq \left(\ell_{f,1} + \frac{\ell_{f,0}\ell_{g,2}}{\mu} + \ell_{g,1} \left(1 + \sqrt{\kappa} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^N \right) \left(\ell_{g,2}\ell_{f,0} + \frac{\ell_{f,1}}{\mu} \right) \right) \text{dist}(y^*(x^k), y^k) \\
 &+ \ell_{g,1} \sqrt{\kappa} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^N \|\hat{v}^0(x^k, y^k) - P_{y^*(x^k) \rightarrow y^k} v^*(x^k)\|.
 \end{aligned}$$

We obtain the desired result by applying Lemma 15 to the above inequality. \blacksquare

The following technical lemma is needed to further bound the right hand side of the inequality in the above lemma.

Lemma 17 *Suppose Assumptions 1, 2, 3, 4 and 5 hold, then RieBO satisfies:*

$$\begin{aligned}
 \text{dist}(y^{k,0}, y^*(x^k))^2 + \|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \hat{v}^0(x^k, y^k)\|^2 \\
 \leq \left(\frac{1}{2} \right)^k \Delta_0 + \Omega \sum_{j=0}^{k-1} \left(\frac{1}{2} \right)^{k-1-j} \|\text{grad}\Phi(x^j)\|^2,
 \end{aligned} \tag{5.14}$$

with the following choice of parameters:

$$\begin{aligned}
 T &\geq \log \left(2 \left(7 + 8\kappa^2 \alpha^2 \Gamma^2 \right) \left(\ell_{g,2}\ell_{f,0} + \frac{\ell_{f,1}}{\mu} \right)^2 / \left(2 \log \left(\frac{1}{1 - 2\mu\tau\beta^2} \right) \right) \right) = \Theta(\kappa), \\
 N &\geq \log \left((4 + 16\kappa^2 \alpha^2 \ell_{g,1}^2) \kappa / \left(2 \log \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right) \right) \right) = \Theta(\sqrt{\kappa}), \\
 \Omega &= \left[2 \left(\frac{\ell_{f,0}\ell_{g,2}}{\mu^2} \sqrt{1 + \kappa^2} + \frac{\ell_{f,1}}{\mu} \right)^2 + 4\kappa^2 \right] \alpha^2, \\
 \Delta_0 &= \text{dist}(y^{0,0}, y^*(x^0))^2 + \|P_{y^*(x^0) \rightarrow y^0} v^*(x^0) - \hat{v}^0(x^0, y^0)\|^2.
 \end{aligned} \tag{5.15}$$

Proof Since $y^{k,0} = y^{k-1,T}$, we have

$$\text{dist}(y^{k,0}, y^*(x^k))^2 \leq 2 \text{dist}(y^{k-1,T}, y^*(x^{k-1}))^2 + 2 \text{dist}(y^*(x^{k-1}), y^*(x^k))^2.$$

Here the first term is again bounded by $(1 - 2\mu\tau\beta^2)^T \text{dist}(y^{k-1,0}, y^*(x^{k-1}))^2$ by Lemma 15, and the second term is bounded by the Lipschitzness of y^* (Lemma 12) and by the update in the following way:

$$\text{dist}(y^*(x^{k-1}), y^*(x^k))^2 \leq \kappa^2 \text{dist}(x^{k-1}, x^k)^2 = \kappa^2 \alpha^2 \|h_{\Phi}^{k-1}\|^2.$$

Thus,

$$\begin{aligned}
 & \text{dist}(y^{k,0}, y^*(x^k))^2 \leq 2 \text{dist}(y^{k-1,T}, y^*(x^{k-1}))^2 + 2 \text{dist}(y^*(x^{k-1}), y^*(x^k))^2 \\
 & \leq 2(1 - 2\mu\tau\beta^2)^T \text{dist}(y^{k-1,0}, y^*(x^{k-1}))^2 + 2\kappa^2\alpha^2 \|h_{\Phi}^{k-1}\|^2 \\
 & \leq 2(1 - 2\mu\tau\beta^2)^T \text{dist}(y^{k-1,0}, y^*(x^{k-1}))^2 + 4\kappa^2\alpha^2 \|h_{\Phi}^{k-1} - \text{grad}\Phi(x^{k-1})\|^2 + 4\kappa^2\alpha^2 \|\text{grad}\Phi(x^{k-1})\|^2 \\
 & \leq \left(2 + 8\kappa^2\alpha^2\Gamma^2\right) (1 - 2\mu\tau\beta^2)^T \text{dist}(y^*(x^{k-1}), y^{k-1,0})^2 \\
 & \quad + 8\kappa^2\alpha^2\ell_{g,1}^2\kappa \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^{2N} \|\hat{v}^0(x^{k-1}, y^{k-1}) - \tilde{v}^{k-1}\|^2 + 4\kappa^2\alpha^2 \|\text{grad}\Phi(x^{k-1})\|^2 \\
 & \leq \left(2 + 8\kappa^2\alpha^2\Gamma^2\right) (1 - 2\mu\tau\beta^2)^T \text{dist}(y^*(x^{k-1}), y^{k-1,0})^2 + 4\kappa^2\alpha^2 \|\text{grad}\Phi(x^{k-1})\|^2 \\
 & \quad + 16\kappa^2\alpha^2\ell_{g,1}^2\kappa \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^{2N} \|\hat{v}^0(x^{k-1}, y^{k-1}) - P_{y^*(x^{k-1}) \rightarrow y^{k-1}} v^*(x^{k-1})\|^2 \\
 & \quad + 16\kappa^2\alpha^2\ell_{g,1}^2\kappa \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^{2N} \|P_{y^*(x^{k-1}) \rightarrow y^{k-1}} v^*(x^{k-1}) - \tilde{v}^{k-1}\|^2,
 \end{aligned}$$

where the third inequality is by Lemma 16. For the last term, by (5.13) we have

$$\|P_{y^*(x^{k-1}) \rightarrow y^k} v^*(x^{k-1}) - \tilde{v}^{k-1}\|^2 \leq \left(\ell_{g,2}\ell_{f,0} + \frac{\ell_{f,1}}{\mu}\right)^2 \text{dist}(y^{k-1}, y^*(x^{k-1}))^2. \quad (5.16)$$

Thus we have

$$\begin{aligned}
 & \text{dist}(y^{k,0}, y^*(x^k))^2 \\
 & \leq \left(2 + 8\kappa^2\alpha^2\Gamma^2\right) (1 - 2\mu\tau\beta^2)^T \text{dist}(y^*(x^{k-1}), y^{k-1,0})^2 + 4\kappa^2\alpha^2 \|\text{grad}\Phi(x^{k-1})\|^2 \\
 & \quad + 16\kappa^2\alpha^2\ell_{g,1}^2\kappa \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^{2N} \|\hat{v}^0(x^{k-1}, y^{k-1}) - P_{y^*(x^{k-1}) \rightarrow y^{k-1}} v^*(x^{k-1})\|^2 \\
 & \quad + 16\kappa^2\alpha^2\ell_{g,1}^2\kappa \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^{2N} \left(\ell_{g,2}\ell_{f,0} + \frac{\ell_{f,1}}{\mu}\right)^2 \text{dist}(y^{k-1}, y^*(x^{k-1}))^2 \\
 & \leq \left[2 + 8\kappa^2\alpha^2\Gamma^2 + 16\kappa^2\alpha^2\ell_{g,1}^2\kappa \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^{2N} \left(\ell_{g,2}\ell_{f,0} + \frac{\ell_{f,1}}{\mu}\right)^2\right] (1 - 2\mu\tau\beta^2)^T \text{dist}(y^*(x^{k-1}), y^{k-1,0})^2 \\
 & \quad + 16\kappa^2\alpha^2\ell_{g,1}^2\kappa \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^{2N} \|\hat{v}^0(x^{k-1}, y^{k-1}) - P_{y^*(x^{k-1}) \rightarrow y^{k-1}} v^*(x^{k-1})\|^2 + 4\kappa^2\alpha^2 \|\text{grad}\Phi(x^{k-1})\|^2.
 \end{aligned}$$

Now we bound $\|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \hat{v}^0(x^k, y^k)\|^2$. We have

$$\begin{aligned}
 & \|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \hat{v}^0(x^k, y^k)\|^2 = \|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - P_{y^{k-1} \rightarrow y^k} \hat{v}^N(x^{k-1}, y^{k-1})\|^2 \\
 & \leq 2\|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - P_{y^*(x^{k-1}) \rightarrow y^k} v^*(x^{k-1})\|^2 \\
 & \quad + 2\|P_{y^*(x^{k-1}) \rightarrow y^k} v^*(x^{k-1}) - P_{y^{k-1} \rightarrow y^k} \hat{v}^N(x^{k-1}, y^{k-1})\|^2 \\
 & \leq 2\|P_{y^*(x^k) \rightarrow y^*(x^{k-1})} v^*(x^k) - v^*(x^{k-1})\|^2 \\
 & \quad + 4\|P_{y^*(x^{k-1}) \rightarrow y^{k-1}} v^*(x^{k-1}) - \tilde{v}^{k-1}\|^2 + 4\|\tilde{v}^{k-1} - \hat{v}^N(x^{k-1}, y^{k-1})\|^2 \\
 & \leq 4\kappa \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2N} \|\tilde{v}^{k-1} - \hat{v}^0(x^{k-1}, y^{k-1})\|^2 \\
 & \quad + 2\|P_{y^*(x^k) \rightarrow y^*(x^{k-1})} v^*(x^k) - v^*(x^{k-1})\|^2 + 4\|P_{y^*(x^{k-1}) \rightarrow y^k} v^*(x^{k-1}) - \tilde{v}^{k-1}\|^2 \\
 & \leq 4\kappa \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2N} \|\tilde{v}^{k-1} - P_{y^*(x^{k-1}) \rightarrow y^{k-1}} v^*(x^{k-1})\|^2 \\
 & \quad + 4\kappa \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2N} \|P_{y^*(x^{k-1}) \rightarrow y^{k-1}} v^*(x^{k-1}) - \hat{v}^0(x^{k-1}, y^{k-1})\|^2 \\
 & \quad + 2\|P_{y^*(x^k) \rightarrow y^*(x^{k-1})} v^*(x^k) - v^*(x^{k-1})\|^2 + 4\|P_{y^*(x^{k-1}) \rightarrow y^{k-1}} v^*(x^{k-1}) - \tilde{v}^{k-1}\|^2 \\
 & = 4\kappa \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2N} \|P_{y^*(x^{k-1}) \rightarrow y^{k-1}} v^*(x^{k-1}) - \hat{v}^0(x^{k-1}, y^{k-1})\|^2 \\
 & \quad + 2\|P_{y^*(x^k) \rightarrow y^*(x^{k-1})} v^*(x^k) - v^*(x^{k-1})\|^2 \\
 & \quad + 4 \left(\kappa \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2N} + 1 \right) \|P_{y^*(x^{k-1}) \rightarrow y^{k-1}} v^*(x^{k-1}) - \tilde{v}^{k-1}\|^2,
 \end{aligned} \tag{5.17}$$

where in the second last inequality we again used Lemma 14. Now we inspect the two terms in the last line above. Note that the last term is bounded in (5.16). For the first term, by (5.7) we have

$$\|P_{y^*(x^k) \rightarrow y^*(x^{k-1})} v^*(x^k) - v^*(x^{k-1})\|^2 \leq \left(\frac{\ell_{f,0} \ell_{g,2}}{\mu^2} \sqrt{1 + \kappa^2} + \frac{\ell_{f,1}}{\mu} \right)^2 \alpha^2 \|\text{grad}\Phi(x^{k-1})\|^2.$$

Now plugging everything back to (5.17) we get

$$\begin{aligned}
 & \|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \hat{v}^0(x^k, y^k)\|^2 \\
 & \leq 4\kappa \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2N} \|P_{y^*(x^{k-1}) \rightarrow y^{k-1}} v^*(x^{k-1}) - \hat{v}^0(x^{k-1}, y^{k-1})\|^2 \\
 & \quad + 2 \left(\frac{\ell_{f,0} \ell_{g,2}}{\mu^2} \sqrt{1 + \kappa^2} + \frac{\ell_{f,1}}{\mu} \right)^2 \alpha^2 \|\text{grad}\Phi(x^{k-1})\|^2 \\
 & \quad + 4 \left(\kappa \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2N} + 1 \right) \left(\ell_{g,2} \ell_{f,0} + \frac{\ell_{f,1}}{\mu} \right)^2 (1 - 2\mu\tau\beta^2)^T \text{dist}(y^{k-1,0}, y^*(x^{k-1}))^2,
 \end{aligned} \tag{5.18}$$

where we also used Lemma 15 in the last inequality. Now summing up the bound for $\text{dist}(y^{k,0}, y^*(x^k))^2$ and $\|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \hat{v}^0(x^k, y^k)\|^2$, we get:

$$\begin{aligned} & \text{dist}(y^{k,0}, y^*(x^k))^2 + \|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \hat{v}^0(x^k, y^k)\|^2 \\ \leq & C_1(1 - 2\mu\tau\beta^2)^T \text{dist}(y^{k-1,0}, y^*(x^{k-1}))^2 \\ & + C_2 \|P_{y^*(x^{k-1}) \rightarrow y^{k-1}} v^*(x^{k-1}) - \hat{v}^0(x^{k-1}, y^{k-1})\|^2 \\ & + \left[2 \left(\frac{\ell_{f,0}\ell_{g,2}}{\mu^2} \sqrt{1 + \kappa^2} + \frac{\ell_{f,1}}{\mu} \right)^2 + 4\kappa^2 \right] \alpha^2 \|\text{grad}\Phi(x^{k-1})\|^2, \end{aligned}$$

with

$$\begin{aligned} C_1 &= \left(6 + 8\kappa^2\alpha^2\Gamma^2 + C_2 \right) \left(\ell_{g,2}\ell_{f,0} + \frac{\ell_{f,1}}{\mu} \right)^2 \\ C_2 &= \left(4 + 16\kappa^2\alpha^2\ell_{g,1}^2 \right) \kappa \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2N}. \end{aligned}$$

Now consider the choice of T and N in the statement of this lemma, we can guarantee that $C_1, C_2 \leq 1/2$, thus

$$\begin{aligned} & \text{dist}(y^{k,0}, y^*(x^k))^2 + \|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \hat{v}^0(x^k, y^k)\|^2 \\ \leq & \frac{1}{2} (\text{dist}(y^{k-1,0}, y^*(x^{k-1}))^2 + \|P_{y^*(x^k) \rightarrow y^k} v^*(x^k) - \hat{v}^0(x^k, y^k)\|^2) + \Omega \|\text{grad}\Phi(x^{k-1})\|^2. \end{aligned}$$

The final result is obtained by taking the telescoping sum of the above inequality. \blacksquare

Combining above Lemma 16 and 17 we get the following lemma.

Lemma 18 *Suppose the parameters are set the same as in Lemma 17, then we have*

$$\|h_{\Phi}^k - \text{grad}\Phi(x^k)\|^2 \leq \delta_{T,N} \left(\frac{1}{2} \right)^k \Delta_0 + \delta_{T,N} \Omega \sum_{j=0}^{k-1} \left(\frac{1}{2} \right)^{k-1-j} \|\text{grad}\Phi(x^j)\|^2, \quad (5.19)$$

where

$$\delta_{T,N} = 2\Gamma^2(1 - 2\mu\tau\beta^2)^T \ell_{g,1}^2 \kappa \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2N}. \quad (5.20)$$

Proof By Lemma 16 and $ab + cd \leq (a + c)(b + d)$ for any positive a, b, c, d , we have

$$\|h_{\Phi}^k - \text{grad}\Phi(x^k)\|^2 \leq \delta_{T,N} \left(\text{dist}(y^*(x^k), y^{k-1})^2 + \|\hat{v}^0(x^k, y^k) - P_{y^*(x^k) \rightarrow y^k} v^*(x^k)\|^2 \right).$$

The proof is completed by applying Lemma 17. \blacksquare

Finally, we are ready to proceed to the proof of Theorem 13.

Proof [Proof of Theorem 13] By Lemma 12, we have

$$\begin{aligned}
 \Phi(x^{k+1}) &\leq \Phi(x^k) + \langle \text{grad}\Phi(x^k), \text{Exp}_{x^k}^{-1}(x^{k+1}) \rangle_{x^k} + \frac{L_\Phi}{2} \text{dist}(x^k, x^{k+1})^2 \\
 &= \Phi(x^k) - \alpha \langle \text{grad}\Phi(x^k), h_\Phi^k \rangle_{x^k} + \frac{L_\Phi \alpha^2}{2} \|h_\Phi^k\|_{x^k}^2 \\
 &\leq \Phi(x^k) - \left(\frac{\alpha}{2} - \alpha^2 L_\Phi\right) \|\text{grad}\Phi(x^k)\|_{x^k}^2 + \left(\frac{\alpha}{2} + \alpha^2 L_\Phi\right) \|\text{grad}\Phi(x^k) - h_\Phi^k\|_{x^k}^2.
 \end{aligned}$$

Now by using Lemma 18, we get

$$\begin{aligned}
 \Phi(x^{k+1}) &\leq \Phi(x^k) - \left(\frac{\alpha}{2} - \alpha^2 L_\Phi\right) \|\text{grad}\Phi(x^k)\|_{x^k}^2 \\
 &\quad + \left(\frac{\alpha}{2} + \alpha^2 L_\Phi\right) \left[\delta_{T,N} \left(\frac{1}{2}\right)^k \Delta_0 + \delta_{T,N} \Omega \sum_{j=0}^{k-1} \left(\frac{1}{2}\right)^{k-1-j} \|\text{grad}\Phi(x^j)\|^2 \right].
 \end{aligned}$$

Now by taking the telescoping sum of the above inequality over k from 0 to $K-1$, we have

$$\begin{aligned}
 \left(\frac{\alpha}{2} - \alpha^2 L_\Phi\right) \sum_{k=0}^{K-1} \|\text{grad}\Phi(x^k)\|^2 &\leq \Phi(x_0) - \inf_{x \in \mathcal{M}} \Phi(x) + \left(\frac{\alpha}{2} + \alpha^2 L_\Phi\right) \delta_{T,N} \Delta_0 \\
 &\quad + \left(\frac{\alpha}{2} + \alpha^2 L_\Phi\right) \delta_{T,N} \Omega \sum_{k=1}^{K-1} \sum_{j=0}^{k-1} \left(\frac{1}{2}\right)^{k-1-j} \|\text{grad}\Phi(x^j)\|^2.
 \end{aligned}$$

By the fact that

$$\sum_{k=1}^{K-1} \sum_{j=0}^{k-1} \left(\frac{1}{2}\right)^{k-1-j} \|\text{grad}\Phi(x^j)\|^2 \leq \sum_{k=0}^{K-1} \frac{1}{2^k} \sum_{k=0}^{K-1} \|\text{grad}\Phi(x^k)\|^2 \leq 2 \sum_{k=0}^{K-1} \|\text{grad}\Phi(x^k)\|^2,$$

we have

$$\left(\frac{\alpha}{2} - \alpha^2 L_\Phi - (\alpha + 2\alpha^2 L_\Phi) \delta_{T,N} \Omega\right) \sum_{k=0}^{K-1} \|\text{grad}\Phi(x^k)\|^2 \leq \Phi(x_0) - \inf_{x \in \mathcal{M}} \Phi(x) + \left(\frac{\alpha}{2} + \alpha^2 L_\Phi\right) \delta_{T,N} \Delta_0.$$

Choosing $N \geq \Theta(\sqrt{\kappa})$ and $D \geq \Theta(\kappa)$ as in Lemma 17, we are able to ensure that

$$\Omega(1 + 2\alpha L_\Phi) \delta_{T,N} \leq \frac{1}{4}, \quad \delta_{T,N} \leq 1.$$

As a result, we get

$$\left(\frac{\alpha}{4} - \alpha^2 L_\Phi\right) \sum_{k=0}^{K-1} \|\text{grad}\Phi(x^k)\|^2 \leq \Phi(x_0) - \inf_{x \in \mathcal{M}} \Phi(x) + \left(\frac{\alpha}{2} + \alpha^2 L_\Phi\right) \Delta_0.$$

Thus, with $\alpha \leq \frac{1}{8L_\Phi}$ we get

$$\frac{1}{K} \sum_{k=0}^{K-1} \|\text{grad}\Phi(x^k)\|^2 \leq \frac{64L_\Phi (\Phi(x_0) - \inf_x \Phi(x)) + 5\Delta_0}{K}.$$

Now we inspect the oracle complexities. To ensure $\frac{1}{K} \sum_{k=0}^{K-1} \|\mathbf{grad}\Phi(x^k)\|^2 \leq \epsilon$, we need $K = \mathcal{O}(\kappa^3/\epsilon)$ where the dependency over κ is due to $L_\Phi = \mathcal{O}(\kappa^3)$ from (5.4), therefore $\text{Gc}(f, \epsilon) = \mathcal{O}(\kappa^3/\epsilon)$. Since in each outer iteration, we need $D = \mathcal{O}(\kappa)$ iterations, we have $\text{Gc}(g, \epsilon) = \mathcal{O}(\kappa^4\epsilon)$. The Jacobian-vector product count is the same as the iteration number K since it is only conducted once for every iteration. The Hessian-vector product is conducted for $N = \mathcal{O}(\sqrt{\kappa})$ times for each iteration. Thus we have the previously described complexities. ■

6. Stochastic Algorithm RieSBO and Its Convergence

In this section, we propose RieSBO (Algorithm 2) for stochastic bilevel manifold optimization (1.2). The algorithm is a generalization of its counterpart in the Euclidean space as in Hong et al. (2023); Chen et al. (2021), where we employ the Neumann series estimation for the hypergradient as in (4.13).

Algorithm 2: Algorithm for Riemannian Stochastic Bilevel Optimization (RieSBO)

```

input :  $K, T, Q$ , stepsize  $\{\alpha_k, \beta_k\}$ , initializations  $x^0 \in \mathcal{M}, y^0 \in \mathcal{N}$ 
for  $k = 0, 1, 2, \dots, K - 1$  do
    Set  $y^{k,0} = y^{k-1}$ ;
    for  $t = 0, \dots, T - 1$  do
        Update  $y^{k,t+1} \leftarrow \text{Exp}_{y^{k,t}}(-\beta_k \tilde{h}_g^{k,t})$  with  $\tilde{h}_g^{k,t} := \mathbf{grad}_y G(x^k, y^{k,t}; \zeta_{k,t})$ ;
    end
    Set  $y^k \leftarrow y^{k,T}$ ;
    Update  $x^{k+1} \leftarrow \text{Exp}_{x^k}(-\alpha_k \tilde{h}_\Phi^k)$ , where  $\tilde{h}_\Phi^k$  is as defined in (4.13);
end
    
```

For the stochastic case, we utilize the following notion of stationarity.

Definition 19 *A random point $x \in \mathcal{M}$ is called an ϵ -stationary point for (1.2) if $\mathbb{E}\|\nabla\Phi(x)\|^2 \leq \epsilon$.*

We now proceed to the convergence analysis for the Riemannian stochastic bilevel optimization (RieSBO, Algorithm 2). For RieSBO, we need the following additional assumption over the mean and variance of the estimators.

Assumption 6 *The stochastic gradients satisfy $\mathbf{grad}F(x, y; \xi) = [\mathbf{grad}_x F(x, y; \xi), \mathbf{grad}_y F(x, y; \xi)]$ and $\mathbf{grad}G(x, y; \zeta) = [\mathbf{grad}_x G(x, y; \zeta), \mathbf{grad}_y G(x, y; \zeta)]$. The second order gradients $\mathbf{grad}_{x,y}^2 G(x, y; \zeta)$, $H_y(G(x, y; \zeta))$ are all unbiased estimators of the corresponding deterministic quantities of f and g . Their variances are all bounded by σ^2 (in tangent space norms and operator norms, respectively for the Riemannian gradient and Riemannian Hessian).*

Note that we do not need to assume the smoothness or strong-convexity of the stochastic functions F and G .

Now we are ready to provide our convergence analysis result for RieSBO (Algorithm 2). This result is an extension of the convergence result in Chen et al. (2021) to the Riemannian setting.

Theorem 20 *Suppose Assumptions 1, 3, 4, 5 and 6 hold. If we take the stepsizes $\alpha_k = \alpha = \frac{1}{\kappa^{5/2}\sqrt{K}}$, $\beta_k = \beta = \min\{\frac{1}{\kappa^{7/4}\sqrt{K}}, \frac{1}{\ell_{g,1}}\}$, also $\eta = 1/\ell_{g,1}$, $Q = \mathcal{O}(\kappa \log K)$ and $T = \mathcal{O}(\kappa^4)$. Also suppose that the random variables for all iterations ζ_k^t , $\zeta_{k,(q)}$, ξ_k are i.i.d. samples, then RieSBO (Algorithm 2) satisfies*

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\text{grad}\Phi(x^k)\|^2] \leq \mathcal{O}\left(\frac{\kappa^{2.5}}{\sqrt{K}}\right).$$

Here the expectation is taken with respect to all the random samples. In order to obtain an ϵ -stationary point, i.e., $\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\text{grad}\Phi(x^k)\|^2] \leq \epsilon$, the oracle complexities needed are given by:

- *Gradients:* $\text{Gc}(f, \epsilon) = \mathcal{O}(\kappa^5 \epsilon^{-2})$, $\text{Gc}(g, \epsilon) = \mathcal{O}(\kappa^9 \epsilon^{-2})$;
- *Jacobian and Hessian-vector products:* $\text{JV}(g, \epsilon) = \mathcal{O}(\kappa^5 \epsilon^{-2})$, $\text{HV}(g, \epsilon) = \tilde{\mathcal{O}}(\kappa^6 \epsilon^{-2})$.

where we hide additional $\log(1/\epsilon)$ factors in $\tilde{\mathcal{O}}$.

To prove this theorem, we need the following lemmas. For simplicity, denote \mathcal{U}_k the σ -algebra generated by all the random samples up to the $(k-1)$ -th iterate, and denote $\bar{h}_\Phi^k := \mathbb{E}[\tilde{h}_\Phi^k | \mathcal{U}_k]$, i.e., the expectation only with respect to the samples of the current iterate.

Lemma 21 *Suppose we estimate the hypergradient \tilde{h}_Φ^k via (4.13) with $\eta \leq \frac{1}{\ell_{g,1}}$, then we have the following bounds.*

$$\mathbb{E}[\|\tilde{h}_\Phi^k - \bar{h}_\Phi^k\|^2 | \mathcal{U}_k] \leq \tilde{\sigma}^2, \quad (6.1)$$

and

$$\|\text{grad}\hat{\Phi}(x^k) - \bar{h}_\Phi^k\|^2 \leq b_k^2, \quad (6.2)$$

where

$$\tilde{\sigma}^2 := 2\sigma^2 + 6(\sigma^2(\sigma^2 + \ell_{f,0}^2) + \ell_{g,1}^2(\sigma^2 + \ell_{f,0}^2) + \ell_{g,1}^2\sigma^2) \max\left\{\frac{1}{\mu^2}, \frac{d_1^2}{\eta^2\mu^2}\right\} = \mathcal{O}(\kappa^2), \quad (6.3)$$

$$b_k := \ell_{f,0} \frac{\ell_{g,1}}{\mu} \left(1 - \frac{\mu}{\ell_{g,1}}\right)^Q,$$

and $\hat{\Phi}(x) = f(x, y^T(x))$ which is the approximate function after T steps of the inner loop.

Further, we have the following bound on the second moment:

$$\mathbb{E}[\|\tilde{h}_\Phi^k\|^2 | \mathcal{U}_k] \leq 2\tilde{\sigma}^2 + 4b_k^2 + 4\ell_{f,0}^2(1 + \kappa)^2 =: \tilde{C}^2 = \mathcal{O}(\kappa^2). \quad (6.4)$$

Proof [Proof of Lemma 21] By the expression (4.1) for $\text{grad}\Phi(x)$ and (4.11) for \tilde{h}_Φ^k , we have:

$$\begin{aligned} \text{grad}\Phi(x^k) &= \text{grad}_x f(x^k, y^*(x^k)) - \text{grad}_{y,x}^2 g(x^k, y^*(x^k))[v^*(x^k)], \\ \text{grad}\hat{\Phi}(x^k) &= \text{grad}_x f(x^k, y^k) - \text{grad}_{y,x}^2 g(x^k, y^k)[\tilde{v}^k], \\ \tilde{h}_\Phi^k &= \text{grad}_x F(x^k, y^k; \xi_k) - \text{grad}_{y,x}^2 G(x^k, y^k; \zeta_{k,(0)})[v_Q^k], \end{aligned}$$

where again $\tilde{v}^k := (H_y(g(x^k, y^{k,T})))^{-1} \mathbf{grad}_y f(x^k, y^{k,T})$.

For (6.1), denote

$$\bar{v}_Q^k = \mathbb{E}[v_Q^k] = \eta \sum_{q=1}^Q (I - \eta H_y(g(x^k, y^k)))^q [\mathbf{grad}_y f(x^k, y^k)].$$

We have

$$\begin{aligned} & \mathbb{E}[\|\tilde{h}_\Phi^k - \bar{h}_\Phi^k\|^2 \mid \mathcal{U}_k] \\ & \leq 2\mathbb{E}[\|\mathbf{grad}_{x,f}(x^k, y^{k,T}) - \mathbf{grad}_x F(x^k, y^{k,T}; \xi_k)\|^2 \mid \mathcal{U}_k] \\ & \quad + 2\mathbb{E}[\|\mathbf{grad}_{y,x}^2 G(x^k, y^k; \zeta_{k,(0)})[v_Q^k] - \mathbf{grad}_{y,x}^2 g(x^k, y^k)[\bar{v}_Q^k]\|^2 \mid \mathcal{U}_k] \\ & \leq 2\sigma^2 + 2\mathbb{E}[\|\mathbf{grad}_{y,x}^2 G(x^k, y^k; \zeta_{k,(0)})[v_Q^k] - \mathbf{grad}_{y,x}^2 g(x^k, y^k)[\bar{v}_Q^k]\|^2 \mid \mathcal{U}_k]. \end{aligned} \tag{6.5}$$

We now inspect the last term above. Denote

$$H^k := \eta Q \prod_{q=1}^{Q'} (I - \eta H_y(G(x^k, y^k; \zeta_{k,(q)}))),$$

which is our estimation of the Riemannian Hessian at the k -th outer iteration, and we have that

$$\begin{aligned} & \mathbf{grad}_{y,x}^2 G(x^k, y^k; \zeta_{k,(0)})[v_Q^k] - \mathbf{grad}_{y,x}^2 g(x^k, y^k)[\bar{v}_Q^k] \\ & = \mathbf{grad}_{y,x}^2 G(x^k, y^k; \zeta_{k,(0)}) \left[H^k [\mathbf{grad}_y F(x^k, y^k; \xi_k)] \right] - \mathbf{grad}_{y,x}^2 g(x^k, y^k) \left[\mathbb{E}[H^k [\mathbf{grad}_y F(x^k, y^k; \xi_k)]] \right] \\ & = \left\{ \mathbf{grad}_{y,x}^2 G(x^k, y^k; \zeta_{k,(0)}) - \mathbf{grad}_{y,x}^2 g(x^k, y^k) \right\} \left[H^k [\mathbf{grad}_y F(x^k, y^k; \xi_k)] \right] \\ & \quad + \mathbf{grad}_{y,x}^2 g(x^k, y^k) \left[\left\{ H^k - \mathbb{E}[H^k] \right\} [\mathbf{grad}_y F(x^k, y^k; \xi_k)] \right] \\ & \quad + \mathbf{grad}_{y,x}^2 g(x^k, y^k) \mathbb{E}[H^k] \left\{ \mathbf{grad}_y F(x^k, y^k; \xi_k) - \mathbf{grad}_y f(x^k, y^k) \right\}. \end{aligned}$$

Since

$$\begin{aligned} & \mathbb{E}[\|\mathbf{grad}_y F(x^k, y^k; \xi_k)\|^2] \\ & = \mathbb{E}[\|\mathbf{grad}_y F(x^k, y^k; \xi_k) - \mathbf{grad}_y f(x^k, y^k)\|^2] + \mathbb{E}[\|\mathbf{grad}_y f(x^k, y^k)\|^2] \leq \sigma^2 + \ell_{f,0}^2, \end{aligned}$$

we have that

$$\begin{aligned} & \mathbb{E}[\|\mathbf{grad}_{y,x}^2 G(x^k, y^k; \zeta_{k,(0)})[v_Q^k] - \mathbf{grad}_{y,x}^2 g(x^k, y^k)[\bar{v}_Q^k]\|^2 \mid \mathcal{U}_k] \\ & \leq 3\sigma^2(\sigma^2 + \ell_{f,0}^2) \mathbb{E}\|H^k\|_{\text{op}}^2 + 3\ell_{g,1}^2(\sigma^2 + \ell_{f,0}^2) \mathbb{E}\|H^k - \mathbb{E}[H^k]\|_{\text{op}}^2 + 3\ell_{g,1}^2 \sigma^2 \|\mathbb{E}[H^k]\|_{\text{op}}^2. \end{aligned}$$

It remains to bound $\mathbb{E}\|H^k\|_{\text{op}}^2$ and $\|\mathbb{E}[H^k]\|_{\text{op}}$. For $\mathbb{E}\|H^k\|_{\text{op}}^2$, using Hong et al. (2023, Lemma 12), we have that

$$\mathbb{E}\|H^k\|_{\text{op}}^2 \leq \frac{d_1}{\eta\mu},$$

where $d_1 > 0$ is some absolute constant. On the other hand $\|\mathbb{E}[H^k]\|_{\text{op}}$ can be easily calculated as (since $\mu\eta < \mu/\ell_{g,1} < 1$)

$$\begin{aligned}\|\mathbb{E}[H^k]\|_{\text{op}} &= \eta \left\| \sum_{q=1}^Q (I - \eta H_y(g(x^k, y^k)))^q \right\|_{\text{op}} \\ &\leq \|H^{-1}\|_{\text{op}} \|I - \eta H_y(g(x^k, y^k))\|_{\text{op}} \leq \frac{1}{\mu}.\end{aligned}$$

Therefore, we finally have

$$\begin{aligned}\mathbb{E}[\|\mathbf{grad}_{y,x}^2 G(x^k, y^k; \zeta_{k,(0)})[v_Q^k] - \mathbf{grad}_{y,x}^2 g(x^k, y^k)[\bar{v}_Q^k]\|^2 \mid \mathcal{U}_k] \\ \leq 3(\sigma^2(\sigma^2 + \ell_{f,0}^2) + \ell_{g,1}^2(\sigma^2 + \ell_{f,0}^2) + \ell_{g,1}^2\sigma^2) \max\left\{\frac{1}{\mu^2}, \frac{d_1^2}{\eta^2\mu^2}\right\}.\end{aligned}$$

Plugging the above equation to (6.5) we get (6.1).

Now for (6.2), since

$$\begin{aligned}\bar{h}_\Phi^k &:= \mathbb{E}\left[\mathbf{grad}_x F(x^k, y^k; \xi_k) - \mathbf{grad}_{y,x}^2 G(x^k, y^k; \zeta_{k,(0)})[v_Q^k]\right] \\ &= \mathbf{grad}_x f(x^k, y^k) - \mathbf{grad}_{y,x}^2 g(x^k, y^k)[\bar{v}_Q^k],\end{aligned}$$

we have

$$\begin{aligned}\|\mathbf{grad}\hat{\Phi}(x^k) - \bar{h}_\Phi^k\|^2 &\leq \|\mathbf{grad}_{y,x}^2 g(x^k, y^k)\|_{\text{op}}^2 \|\tilde{v}^k - \bar{v}_Q^k\|^2 \leq \ell_{g,1}^2 \|\tilde{v}^k - \bar{v}_Q^k\|^2 \\ &= \ell_{g,1}^2 \|(H_y(g(x^k, y^k)))^{-1}[\mathbf{grad}_y f(x^k, y^k)] - \eta \sum_{q=1}^Q (I - \eta H_y(g(x^k, y^k)))^q [\mathbf{grad}_y f(x^k, y^k)]\|^2 \\ &\leq \ell_{g,1}^2 \ell_{f,0}^2 \|(H_y(g(x^k, y^k)))^{-1} - \eta \sum_{q=1}^Q (I - \eta H_y(g(x^k, y^k)))^q\|_{\text{op}}^2 \\ &\leq \ell_{f,0}^2 \frac{\ell_{g,1}^2}{\mu^2} \left(1 - \frac{\mu}{\ell_{g,1}}\right)^{2Q} = b_k^2,\end{aligned}$$

where the last line is by Ghadimi and Wang (2018, Lemma 3.2). Note that we take $\eta \leq \frac{1}{\ell_{g,1}}$ so that the Neumann sequence converges.

Now for the moment $\mathbb{E}[\|\tilde{h}_\Phi^k\|^2 \mid \mathcal{U}_k]$, we have

$$\begin{aligned}\mathbb{E}[\|\tilde{h}_\Phi^k\|^2 \mid \mathcal{U}_k] &\leq 2\mathbb{E}[\|\tilde{h}_\Phi^k - \bar{h}_\Phi^k\|^2 \mid \mathcal{U}_k] + 4\|\bar{h}_\Phi^k - \mathbf{grad}\hat{\Phi}(x^k)\|^2 + 4\|\mathbf{grad}\hat{\Phi}(x^k)\|^2 \\ &\leq 2\tilde{\sigma}^2 + 4b_k^2 + 4\|\mathbf{grad}\hat{\Phi}(x^k)\|^2.\end{aligned}$$

Since

$$\begin{aligned}\|\mathbf{grad}\hat{\Phi}(x^k)\| &= \|\mathbf{grad}_x f(x^k, y^k) - \mathbf{grad}_{y,x}^2 g(x^k, y^k)[\tilde{v}^k]\| \\ &\leq \|\mathbf{grad}_x f(x^k, y^k)\| + \|\mathbf{grad}_{y,x}^2 g(x^k, y^k)\|_{\text{op}} \|\tilde{v}^k\| \leq \ell_{f,0} + \ell_{g,1} \frac{\ell_{f,0}}{\mu} = \ell_{f,0}(1 + \kappa),\end{aligned}$$

we have

$$\mathbb{E}[\|\tilde{h}_\Phi^k\|^2 \mid \mathcal{U}_k] \leq 2\tilde{\sigma}^2 + 4b_k^2 + 4\ell_{f,0}^2(1 + \kappa)^2.$$

This completes the proof. ■

The following lemma quantifies the convergence of the lower level update in RieSBO (Algorithm 2).

Lemma 22 *Suppose we have the sequence $\{y^{k,t}\}$ by RieSBO with stepsize $\beta_k = \beta \leq \frac{1}{\ell_{g,1}}$, then the following inequalities hold:*

$$\mathbb{E} \text{dist}(y^{k,T}, y^*(x^k))^2 \leq (1 - 2\mu\tau\beta^2)^T \text{dist}(y^{k,0}, y^*(x^k))^2 + \tau\beta^2\sigma^2T, \quad (6.6)$$

and

$$\begin{aligned} & \mathbb{E}[\text{dist}(y^{k,T}, y^*(x^{k+1}))^2] \\ & \leq 2(1 - 2\mu\tau\beta^2)^T \text{dist}(y^{k,0}, y^*(x^k))^2 + 2\tau\beta^2\sigma^2T + 4\tau\kappa^2\alpha^2\|\bar{h}_\Phi^k\|_{x^k}^2 + 4\tau\kappa^2\alpha^2\tilde{\sigma}^2. \end{aligned} \quad (6.7)$$

Proof [Proof of Lemma 22] For simplicity, all the expectations are conditioned on \mathcal{U}_k in this proof.

First we have by Zhang and Sra (2016, Corollary 8) that

$$\begin{aligned} & \mathbb{E}_{\zeta_{k,t}} \text{dist}(y^{k,t+1}, y^*(x^k))^2 \\ & \leq \text{dist}(y^{k,t}, y^*(x^k))^2 + 2\beta\langle \text{grad}g(x^k, y^k), \text{Exp}_{y_{k,t}}(y^*(x^k)) \rangle + \tau\beta^2\mathbb{E}_{\zeta_{k,t}}\|\tilde{h}_g^{k,t}\|^2 \\ & \leq \text{dist}(y^{k,t}, y^*(x^k))^2 + 2\beta\langle \text{grad}g(x^k, y^k), \text{Exp}_{y_{k,t}}(y^*(x^k)) \rangle + \tau\beta^2\|\text{grad}g(x^k, y^k)\|^2 + \tau\beta^2\sigma^2 \\ & \leq (1 - 2\mu\tau\beta^2) \text{dist}(y^{k,t}, y^*(x^k))^2 + \tau\beta^2\sigma^2, \end{aligned}$$

where in the last line we used the same trick as the proof of Lemma 15. Note that in the above formulas the expectation is only taken with respect to the random variables in $\tilde{h}_g^{k,t}$, i.e., $\zeta_{k,t}$. Repeating this for T times yields (6.6). Now for the second inequality (6.7), we have

$$\begin{aligned} & \mathbb{E}[\text{dist}(y^{k,T}, y^*(x^{k+1}))^2] \\ & \leq 2\mathbb{E}[\text{dist}(y^{k,T}, y^*(x^k))^2] + 2\mathbb{E} \text{dist}(y^*(x^k), y^*(x^{k+1}))^2 \\ & \leq 2(1 - 2\mu\tau\beta^2)^T \text{dist}(y^{k,0}, y^*(x^k))^2 + 2\tau\beta^2\sigma^2T + 2\tau\kappa^2\mathbb{E} \text{dist}(x^k, x^{k+1})^2, \end{aligned} \quad (6.8)$$

where the last inequality is by (6.6) and Lemma 12. For $\mathbb{E}d(x^{k+1}, x^k)^2$ we have the bound:

$$\begin{aligned} & \mathbb{E}d(x^{k+1}, x^k)^2 = \alpha^2\mathbb{E}\|\tilde{h}_\Phi^k\|_{x^k}^2 \\ & = \alpha^2\mathbb{E}\|\tilde{h}_\Phi^k - \bar{h}_\Phi^k + \bar{h}_\Phi^k\|_{x^k}^2 \leq 2\alpha^2(\|\bar{h}_\Phi^k\|_{x^k}^2 + \tilde{\sigma}^2), \end{aligned}$$

which completes the proof. ■

Now we turn to the proof of Theorem 20.

Proof [Proof of Theorem 20] Denote $V_k := \Phi(x^k) + \kappa \text{dist}(y^{k-1,T}, y^*(x^k))^2$. By Lemma 12 and Lemma 21, we have

$$\begin{aligned}
 \mathbb{E}[\Phi(x^{k+1}) \mid \mathcal{U}_k] &\leq \Phi(x^k) + \mathbb{E}[\langle \text{grad}\Phi(x^k), \text{Exp}_{x^k}^{-1}(x^{k+1}) \rangle_{x^k} \mid \mathcal{U}_k] + \frac{L_\Phi}{2} \mathbb{E}[\text{dist}(x^k, x^{k+1})^2 \mid \mathcal{U}_k] \\
 &= \Phi(x^k) - \alpha \mathbb{E}[\langle \text{grad}\Phi(x^k), \tilde{h}_\Phi^k \rangle_{x^k} \mid \mathcal{U}_k] + \frac{L_\Phi \alpha^2}{2} \|\tilde{h}_\Phi^k\|_{x^k}^2 \\
 &= \Phi(x^k) - \frac{\alpha}{2} \mathbb{E}[\|\text{grad}\Phi(x^k)\|^2 \mid \mathcal{U}_k] - \left(\frac{\alpha}{2} - \frac{\alpha^2 L_\Phi}{2}\right) \|\bar{h}_\Phi^k\|^2 + \frac{\alpha}{2} \|\text{grad}\Phi(x^k) - \bar{h}_\Phi^k\|^2 \\
 &\quad + \frac{\alpha^2 L_\Phi}{2} \mathbb{E}[\|\tilde{h}_\Phi^k - \bar{h}_\Phi^k\|^2 \mid \mathcal{U}_k] \\
 &\leq \Phi(x^k) - \frac{\alpha}{2} \mathbb{E}[\|\text{grad}\Phi(x^k)\|^2 \mid \mathcal{U}_k] - \left(\frac{\alpha}{2} - \frac{\alpha^2 L_\Phi}{2}\right) \|\bar{h}_\Phi^k\|^2 + \frac{\alpha}{2} \|\text{grad}\Phi(x^k) - \bar{h}_\Phi^k\|^2 + \frac{\alpha^2 L_\Phi}{2} \tilde{\sigma}^2.
 \end{aligned}$$

Now we decompose the bias term $\|\text{grad}\Phi(x^k) - \bar{h}_\Phi^k\|$ as:

$$\begin{aligned}
 \|\text{grad}\Phi(x^k) - \bar{h}_\Phi^k\|^2 &= 2\|\text{grad}\Phi(x^k) - \text{grad}\hat{\Phi}(x^k)\|^2 + 2\|\text{grad}\hat{\Phi}(x^k) - \bar{h}_\Phi^k\|^2 \\
 &\leq 2\Gamma_0^2 \text{dist}(y^{k,T}, y^*(x^k))^2 + 2b_k^2,
 \end{aligned} \tag{6.9}$$

where we use a similar process as the proof of Lemma 17 to bound $\|\text{grad}\Phi(x^k) - \text{grad}\hat{\Phi}(x^k)\|$ and $\Gamma_0 = \ell_{f,1} + \frac{\ell_{f,0}\ell_{g,2}}{\mu} + \ell_{g,1}(\ell_{g,2}\ell_{f,0} + \frac{\ell_{f,1}}{\mu}) = \mathcal{O}(\kappa)$. Thus we have

$$\begin{aligned}
 \mathbb{E}[\Phi(x^{k+1}) \mid \mathcal{U}_k] &\leq \Phi(x^k) - \frac{\alpha}{2} \mathbb{E}[\|\text{grad}\Phi(x^k)\|^2 \mid \mathcal{U}_k] - \left(\frac{\alpha}{2} - \frac{\alpha^2 L_\Phi}{2}\right) \|\bar{h}_\Phi^k\|^2 \\
 &\quad + \alpha\Gamma_0^2 \text{dist}(y^{k,T}, y^*(x^k))^2 + \alpha b_k^2 + \frac{\alpha^2 L_\Phi}{2} \tilde{\sigma}^2.
 \end{aligned} \tag{6.10}$$

Now we have

$$\begin{aligned}
 \mathbb{E}[V_{k+1}] - \mathbb{E}[V_k] &= \mathbb{E}[\Phi(x^{k+1})] - \mathbb{E}[\Phi(x^k)] + \kappa \mathbb{E} \text{dist}(y^{k,T}, y^*(x^{k+1}))^2 - \kappa \mathbb{E} \text{dist}(y^{k-1,T}, y^*(x^k))^2 \\
 &\leq -\frac{\alpha}{2} \mathbb{E}[\|\text{grad}\Phi(x^k)\|^2 \mid \mathcal{U}_k] - \left(\frac{\alpha}{2} - \frac{\alpha^2 L_\Phi}{2}\right) \mathbb{E} \|\bar{h}_\Phi^k\|^2 + \alpha b_k^2 + \frac{\alpha^2 L_\Phi}{2} \tilde{\sigma}^2 \\
 &\quad + \kappa \mathbb{E} \text{dist}(y^{k,T}, y^*(x^{k+1}))^2 - \kappa \mathbb{E} \text{dist}(y^{k-1,T}, y^*(x^k))^2 + \alpha\Gamma_0^2 \mathbb{E} \text{dist}(y^{k,T}, y^*(x^k))^2 \\
 &\leq -\frac{\alpha}{2} \mathbb{E}[\|\text{grad}\Phi(x^k)\|^2 \mid \mathcal{U}_k] - \left(\frac{\alpha}{2} - \frac{\alpha^2 L_\Phi}{2}\right) \mathbb{E} \|\bar{h}_\Phi^k\|^2 + \alpha b_k^2 + \frac{\alpha^2 L_\Phi}{2} \tilde{\sigma}^2 \\
 &\quad + \kappa \mathbb{E} \left(2(1 - 2\mu\tau\beta^2)^T \text{dist}(y^{k,0}, y^*(x^k))^2 + 2\tau\beta^2\sigma^2 T + 4\tau\kappa^2\alpha^2 \|\bar{h}_\Phi^k\|_{x^k}^2 + 4\tau\kappa^2\alpha^2 \tilde{\sigma}^2 \right) \\
 &\quad - \kappa \mathbb{E} \text{dist}(y^{k-1,T}, y^*(x^k))^2 + \alpha\Gamma_0^2 \left((1 - 2\mu\tau\beta^2)^T \mathbb{E} \text{dist}(y^{k,0}, y^*(x^k))^2 + \tau\beta^2\sigma^2 T \right) \\
 &= -\frac{\alpha}{2} \mathbb{E}[\|\text{grad}\Phi(x^k)\|^2 \mid \mathcal{U}_k] - \left(\frac{\alpha}{2} - \frac{\alpha^2 L_\Phi}{2} - 4\tau\kappa^3\alpha^2 \right) \mathbb{E} \|\bar{h}_\Phi^k\|^2 \\
 &\quad + \left((2\kappa + \alpha\Gamma_0^2)(1 - 2\mu\tau\beta^2)^T - \kappa \right) \mathbb{E} \text{dist}(y^{k-1,T}, y^*(x^k))^2 \\
 &\quad + \left(2\kappa + \alpha\Gamma_0^2 \right) \tau\beta^2\sigma^2 T + \alpha b_k^2 + \left(\frac{L_\Phi}{2} + 4\tau\kappa^2 \right) \alpha^2 \tilde{\sigma}^2,
 \end{aligned}$$

where the first inequality is by (6.10) and the second inequality is by Lemma 22, as well as the fact that $y^{k+1,0} = y^{k,T}$. To make the coefficients negative, notice that by taking

$$\alpha \leq \frac{1}{L_\Phi + 8\tau\kappa^2}$$

$$T \geq \log\left(\frac{1}{1 - 2\mu\tau\beta^2}\right) / \log\left(\frac{2\kappa + \alpha\Gamma_0^2}{\kappa}\right),$$

we can guarantee

$$\frac{\alpha}{2} - \frac{\alpha^2 L_\Phi}{2} - 4\tau\kappa^3 \alpha^2 \geq 0$$

$$(2\kappa + \alpha\Gamma_0^2)(1 - 2\mu\tau\beta^2)^T - \kappa \leq 0.$$

Therefore, we have

$$\begin{aligned} \mathbb{E}[V_{k+1}] - \mathbb{E}[V_k] &\leq -\frac{\alpha}{2} \mathbb{E}[\|\text{grad}\Phi(x^k)\|^2 \mid \mathcal{U}_k] \\ &\quad + \left(2\kappa + \alpha\Gamma_0^2\right)\tau\beta^2\sigma^2T + \alpha b_k^2 + \left(\frac{L_\Phi}{2} + 4\tau\kappa^2\right)\alpha^2\tilde{\sigma}^2. \end{aligned} \quad (6.11)$$

Note that here we do not need an increasing T .

Now taking the telescoping sum of the above inequality for $k = 0, \dots, K-1$, we get

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\text{grad}\Phi(x^k)\|^2] \leq \frac{2V_0}{\alpha K} + \frac{2}{K} \sum_{k=0}^{K-1} b_k^2 + \left(\frac{4\kappa}{\alpha} + 2\Gamma_0^2\right)\tau\beta^2\sigma^2T + (L_\Phi + 8\tau\kappa^2)\alpha\tilde{\sigma}^2.$$

Now since $b_k^2 = \ell_{f,0}^2 \kappa^2 (1 - \frac{1}{\kappa})^{2Q}$, the term $\frac{2}{K} \sum_{k=0}^{K-1} b_k^2 = \mathcal{O}(\frac{1}{\sqrt{K}})$ if $Q = \mathcal{O}(\kappa \log(K))$, following the inequality $(1-x)^n \leq e^{-nx}$. If we also select $\alpha_k = \alpha = \frac{1}{\kappa^{5/2}\sqrt{K}}$, $\beta_k = \beta = \min\{\frac{1}{\kappa^{7/4}\sqrt{K}}, \frac{1}{\ell_{g,1}}\}$ and $T = \mathcal{O}(\kappa^4)$, we are able to get:

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\text{grad}\Phi(x^k)\|^2] \leq \mathcal{O}\left(\frac{\kappa^{2.5}}{\sqrt{K}}\right).$$

Now we inspect the oracle complexities. To ensure $\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\text{grad}\Phi(x^k)\|^2] \leq \epsilon$, we need $K = \mathcal{O}(\kappa^5 \epsilon^{-2})$, thus $\text{Gc}(F, \epsilon) = \mathcal{O}(\kappa^5 \epsilon^{-2})$; Also $\text{Gc}(G, \epsilon) = KT = \mathcal{O}(\kappa^9 \epsilon^{-2})$. \blacksquare

Remark 23 *Note that the trick for estimating the Hessian-vector product (4.12) can also be applied to the deterministic case, without using conjugate gradient method, leading to an easier implementation. We just need to replace the stochastic functions in (4.12) by their deterministic versions. In the experiments, we always use (4.12) in this way instead of solving (4.8) which uses N -step conjugate gradient method, while still achieving reasonable results numerically.*

7. Numerical experiments

In this section we numerically verify the effectiveness of the proposed RieBO and RieSBO. Our code is publicly available at https://github.com/JasonJiaxiangLi/Manifold_bilevel.

7.1 Numerical experiments on robust optimization on manifolds

Consider again the robust optimization on manifolds:

$$\begin{aligned} \min_{p \in \Delta_n} \quad & \lambda \left\| p - \frac{\mathbf{1}}{n} \right\|^2 - \sum_{i=1}^n p_i \ell(y; \xi_i) \\ \text{s.t.} \quad & y \in \operatorname{argmin}_{y \in \mathcal{N}} \sum_{i=1}^n p_i \ell(y; \xi_i). \end{aligned} \tag{7.1}$$

It is worth noticing that having a constraint set in the upper level problem is not covered in our theoretical analysis due to the fact that existing constrained Riemannian optimization techniques such as (stochastic) Riemannian Frank-Wolfe (see Weber and Sra (2022, 2023)) require a mini-batch sampling technique, i.e., using (4.13) multiple times and taking the average of these estimators to estimate $\operatorname{grad}\Phi(x^k)$ to reduce the variance, which is not desirable in practice. Instead, we point out that since the upper level in (7.1) is a constrained optimization in a Euclidean space, one could utilize the analysis in Hong et al. (2023) to achieve a similar convergence result as the unconstrained case in (1.1). Therefore we simply add a projection step for the upper level update, and we still observed reasonable convergence results. We present the algorithm we use for the numerical experiments in Algorithm 3, which is a direct adaptation of Algorithm 1 to (7.1). Note that here the variables in the upper and lower level problems are respectively denoted by y and x , which is different from the previous algorithms. It is also worth noticing that the convergence criteria are altered due to the existence of the projection step: in Algorithm 3, we simply measure the norm of the quantity

$$\mathcal{G}^k := \frac{1}{\alpha_k} (p^k - p^{k+1}) = \frac{1}{\alpha_k} (p^k - \operatorname{proj}_{\Delta_n}(p^k - \alpha_k h_{\Phi}^k)),$$

which we refer to as the approximate gradient mapping. This quantity can be used for approximately measuring the stationarity since if we do not have a constraint,

$$\mathcal{G}^k = \frac{1}{\alpha_k} (p^k - p^{k+1}) = h_{\Phi}^k \approx \operatorname{grad}\Phi(y^k),$$

based on Lemma 21.

We test our proposed algorithms on two concrete examples that lie in this scope: the robust Karcher mean problem and the robust covariance matrix estimation problems. In both experiments we only consider the deterministic function to test the efficacy of the proposed algorithm framework. In both experiments, we use RieBO (Algorithm 1) while utilizing the trick in Remark 23 to estimate the Hessian-vector products.

Algorithm 3: Bilevel algorithm for robust manifold optimization problem (7.1)

input : K, T, N (steps for conjugate gradient), stepsize $\{\alpha_k, \beta_k\}$, initializations
 $p^0 \in \Delta_n, y^0 \in \mathcal{N}$
for $k = 0, 1, 2, \dots, K - 1$ **do**
 Set $y^{k,0} = y^{k-1}$;
 for $t = 0, \dots, T - 1$ **do**
 Update $y^{k,t+1} \leftarrow \text{Exp}_{y^{k,t}}(-\beta_k h_g^{k,t})$ with $h_g^{k,t} := \text{grad}_y g(p^k, y^{k,t})$;
 end
 Set $y^k \leftarrow y^{k,T}$;
 Update $p^{k+1} \leftarrow \text{proj}_{\Delta_n}(p^k - \alpha_k h_{\Phi}^k)$ as in (4.12), in the view of Remark 23;
end

7.1.1 ROBUST KARCHER MEAN PROBLEM

For the robust Karcher mean problem, one seeks to solve

$$\begin{aligned}
 \min_{p \in \Delta_n} \quad & \left\| p - \frac{\mathbf{1}}{n} \right\|^2 - \sum_{i=1}^n p_i \text{dist}(S, A_i)^2 \\
 \text{s.t. } \quad & S \in \underset{S \in \mathbb{S}_{++}^d}{\text{argmin}} \sum_{i=1}^n y_i \text{dist}(S, A_i)^2,
 \end{aligned} \tag{7.2}$$

where A_i 's are the symmetric positive definite data matrices, and

$$\text{dist}(A, B) := \|\log(A^{-1/2} B A^{-1/2})\|_F$$

is the geodesic distance of two positive definite matrices (see Bhatia (2009, Chapter 6)). The squared geodesic distance guarantees the geodesic strong convexity of the lower level problem (see Zhang and Sra (2016)), which further ensures that the bilevel problem (7.2) is well-defined. For the function $h(S) := \text{dist}(S, A)^2$, we have the Euclidean and Riemannian gradients as (see Ferreira et al. (2019), Bhatia (2009, Chapter 6)):

$$\begin{aligned}
 \nabla_S h(S) &= S^{-1/2} \log(S^{1/2} A^{-1} S^{1/2}) S^{-1/2}, \\
 \text{grad}_S h(S) &= S \nabla_S h(S) S = S^{1/2} \log(S^{1/2} A^{-1} S^{1/2}) S^{1/2}.
 \end{aligned}$$

The Euclidean and Riemannian Hessian of $h(S) := \text{dist}(S, A)^2$ are less straightforward to calculate, and to the best of our knowledge, they do not exist in the literature. Here we propose an implementable way to calculate it: first notice that (see Bhatia (2009, Chapter 6))

$$\nabla_S h(S) = S^{-1/2} \log(S^{1/2} A^{-1} S^{1/2}) S^{-1/2} = S^{-1} A^{1/2} \log(A^{-1/2} S A^{-1/2}) A^{-1/2}.$$

For any symmetric matrix V , we have

$$\langle \nabla_S h(S), V \rangle = \text{tr}(V S^{-1} A^{1/2} \log(A^{-1/2} S A^{-1/2}) A^{-1/2}).$$

To take the derivative of this, notice that S appears twice in $\langle \nabla_S h(S), V \rangle$. Denote $\tilde{h}(S_1, S_2) = \text{tr}(VS_1^{-1}A^{1/2}\log(A^{-1/2}S_2A^{-1/2})A^{-1/2})$, we know that (see Petersen et al. (2008))

$$\frac{\partial \tilde{h}}{\partial S_1} = -S_1^{-1}VA^{-1/2}\log(A^{-1/2}S_2A^{-1/2})A^{1/2}S_1^{-1}.$$

It remains to calculate $\partial \tilde{h}/\partial S_2$, which takes a form of $l(S) := \text{tr}(C\log(PSQ))$. Denote $Y = PSQ$ and $L = \log(Y)$, we have

$$\begin{bmatrix} L & dL \\ 0 & L \end{bmatrix} = \log \left(\begin{bmatrix} Y & dY \\ 0 & Y \end{bmatrix} \right) = \log \left(\begin{bmatrix} P & 0 \\ 0 & P \end{bmatrix} \begin{bmatrix} S & dS \\ 0 & S \end{bmatrix} \begin{bmatrix} Q & 0 \\ 0 & Q \end{bmatrix} \right).$$

Therefore, we get

$$dl = \left\langle \begin{bmatrix} 0 & C \\ 0 & 0 \end{bmatrix}, \log \left(\begin{bmatrix} P & 0 \\ 0 & P \end{bmatrix} \begin{bmatrix} S & dS \\ 0 & S \end{bmatrix} \begin{bmatrix} Q & 0 \\ 0 & Q \end{bmatrix} \right) \right\rangle,$$

where the inner product is simply the Euclidean inner product. We can plug dS as standard Euclidean basis to obtain an representation of dl/dS , which will take $\mathcal{O}(d^2)$ number of times to cover all the entries. Nevertheless, this provides an implementable way to calculate the Euclidean and Riemannian Hessian.

To summarize, the Euclidean and Riemannian Hessian of h can be calculated as follows.

$$\begin{aligned} \nabla_S^2 h(S)[V] &= -S^{-1}VA^{-1/2}\log(A^{-1/2}SA^{-1/2})A^{1/2}S^{-1} + L, \\ H_S(h(S))[V] &= S\nabla_S^2 h(S)[V]S + \text{sym}(S\nabla_S h(S)V), \end{aligned} \tag{7.3}$$

where each entry of matrix L is calculated as follows:

$$L_{i,j} = \left\langle \begin{bmatrix} 0 & C \\ 0 & 0 \end{bmatrix}, \log \left(\begin{bmatrix} P & 0 \\ 0 & P \end{bmatrix} \begin{bmatrix} S & E_{i,j} \\ 0 & S \end{bmatrix} \begin{bmatrix} Q & 0 \\ 0 & Q \end{bmatrix} \right) \right\rangle,$$

Here the (i, j) -th entry of $E_{i,j} \in \mathbb{R}^{d \times d}$ is one, and all other entries are zeros. Moreover,

$$\begin{aligned} P &= A^{-1/2}, \\ Q &= A^{-1/2}, \\ C &= A^{-1/2}VS^{-1}A^{1/2}. \end{aligned}$$

In the experiment, we test RieBO (Algorithm 1) with $d \in \{10, 20\}$ and $n = 5$. We repeat each dimension settings for 5 times and plot the average. The algorithm is terminated with $K = 200$ rounds of outer iterations, and the inner iteration is also taken to be $T = 200$ (the value which we observe a good inner iteration convergence). We take $\alpha_k = 10^{-2}$ and $\beta_k = 10^{-1}$. Figure 1 shows the results of the robust Karcher mean problem (7.2). It can be seen from Figure 1 that Algorithm 3 can efficiently decrease both the function values and the norm of gradient mappings. We point out here that the computation of the Riemannian Hessian is time consuming by (7.3) (which is also the reason why we cannot try larger dimensions), yet we remind the reader that this is currently the only formula for calculating it.

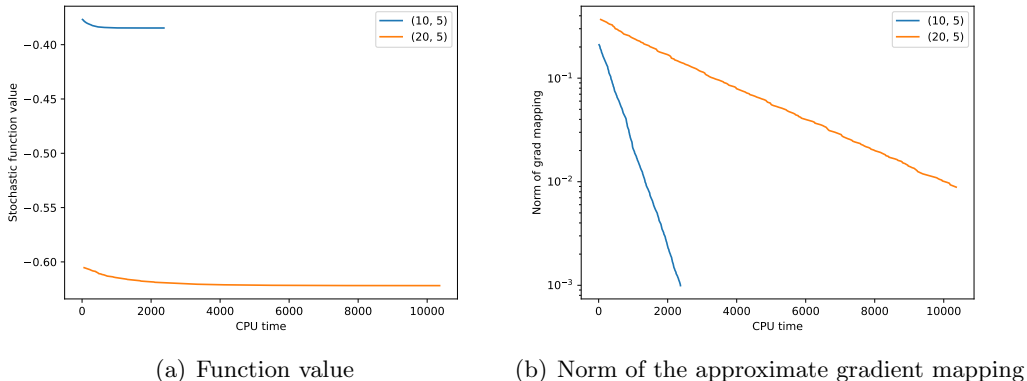


Figure 1: The convergence curve of applying Algorithm 3 to the robust Karcher mean problem (7.2). The CPU time is in seconds.

7.1.2 ROBUST MAXIMUM LIKELIHOOD ESTIMATION

For the robust maximum likelihood estimation of the covariance matrix, one seeks to solve:

$$\begin{aligned}
 \min_{p \in \Delta_n} & \left\| p - \frac{\mathbf{1}}{n} \right\|^2 - \sum_{i=1}^n p_i \mathcal{L}(S; x_i) \\
 \text{s.t. } & S \in \operatorname{argmin}_{S \in \mathbb{S}_{++}^d} \sum_{i=1}^n p_i \mathcal{L}(S; x_i),
 \end{aligned} \tag{7.4}$$

where $\mathcal{L}(S; x)$ is the log likelihood of the Gaussian distribution, namely

$$\mathcal{L}(S; \mathcal{D}) := \frac{1}{2} \log \det(S) + \frac{x^\top S^{-1} x}{2}. \tag{7.5}$$

Note that this lower level problem is geodesically strictly convex (see Sra and Hosseini (2015)), and thus has a unique solution. The calculations of the Riemannian gradient, Hessian-vector product and cross-derivatives all have closed form solutions (following Petersen et al. (2008)).

In the experiment, we test our algorithm with $d \in \{10, 30, 50\}$ and $n = 100$. We repeat each dimension settings for 5 times and plot the average. The algorithm is terminated with $K = 1000$ rounds of outer iterations, and the inner iteration is still taken to be $T = 200$ (again a value which we observe a good inner iteration convergence). We take $\alpha_k = 10^{-2}$ and $\beta_k = 10^{-1}$. Figure 2 shows the results when applying RieBO to the above robust MLE problem with different choices of dimensions. It can be seen from Figure 2 that Algorithm 3 can efficiently decrease both the function values and the norm of gradient mappings. Also, here we are able to test and present the results for a much larger dimension due to much faster calculations of Riemannian gradients, Hessian-vector products and cross-derivatives.

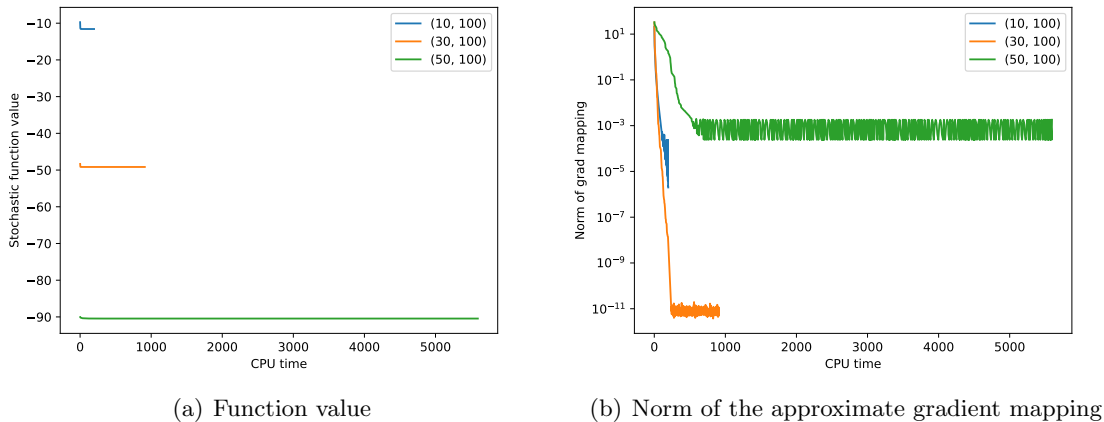


Figure 2: The convergence curve of Algorithm 3 applying to the robust covariance matrix maximum likelihood estimation problem (7.4) with different choice of (d, n) . The CPU time is in seconds.

7.2 Numerical experiments on Riemannian meta learning

We conduct experiment on the Riemannian meta learning problem (2.3) with \mathcal{M} being the Grassmannian manifold $\text{Gr}(n, p)$, which is the manifold of p dimensional subspaces in n dimensional Euclidean space. Such a manifold can be interpreted as the quotient manifold of Stiefel manifold over the orthogonal group, namely $\text{Gr}(n, p) = \text{St}(n, p)/O(p)$ where each element $[X]$ in $\text{Gr}(n, p)$ is an equivalent class $[X] = \{XQ | Q \in O(p)\}$ for an $X \in \text{St}(n, p)$. Usually people use this element $X \in \text{St}(n, p)$ to represent the Grassmannian and employ the projection $\pi : \text{St}(n, p) \rightarrow \text{Gr}(n, p) : X \mapsto \pi(X) \triangleq [X] = \{XQ : Q \in O(p)\}$ to correspond it to the Grassmannian. For Grassmannian $\text{Gr}(n, p)$, the retraction operator would be different from that of Stiefel manifold. To keep the conciseness of this work, we refer to Boumal (2023, Chapter 9) for more details on Grassmannian.

We employ our RieSBO (Algorithm 2) and compare it with a naive projection-based stochastic bilevel algorithm (Denoted *Projected SBO* in our plots), which is basically Algorithm 1 in Ji et al. (2021), with an extra projection onto the Grassmannian $\text{Gr}(n, p)$ at the end of each update step.

Following Li et al. (2020a); Han et al. (2024), we consider 5-ways 5-shots meta learning over the MiniImageNet dataset with four-layer CNN and with the kernels setting to be on the Grassmannian manifold, and test the accuracy over 200 tasks. The result is included in Figure 3, where the RieSBO algorithm is showing better performance in terms of both the training loss and the test accuracy.

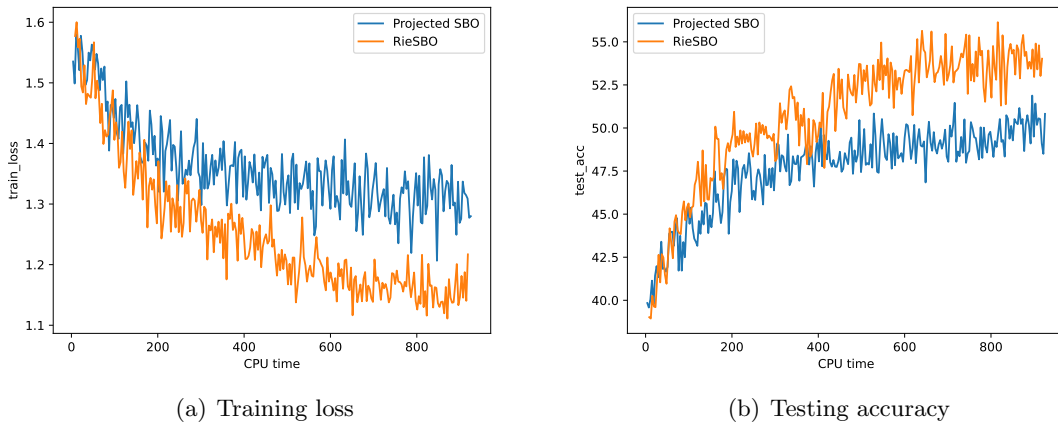


Figure 3: The convergence curve of Algorithm 2 on the meta learning problem.

8. Conclusion

We introduced the Riemannian bilevel optimization, a generalization of the traditional Euclidean bilevel optimization. We show that the Riemannian counterparts of Euclidean algorithms in Chen et al. (2021); Ji et al. (2021) can achieve the same rate of convergence.

Our work raises several open questions. The first is how we can make the convergence independent of the sectional curvature of the manifold, similar to the results in Cai et al. (2023). It is also worth exploring the last iterate convergence of Riemannian bilevel problem. Last, it still needs investigation to see if there are efficient algorithms that can overcome the difficulty we mentioned in the numerical experiment part to efficiently calculate the Riemannian Hessian-vector product thus enabling large-scale implementation of algorithms for solving the Riemannian bilevel optimization problems.

Acknowledgments

JL is supported by NSF grant ECCS-2426064. SM is supported in part by ONR grant N00014-24-1-2705, NSF grants DMS-2243650, CCF-2308597, CCF-2311275 and ECCS-2326591, and a startup fund from Rice University.

References

- P-A Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton University Press, 2008.
- Tamir Bendory, Yonina C Eldar, and Nicolas Boumal. Non-convex phase retrieval from STFT measurements. *IEEE Transactions on Information Theory*, 64(1):467–484, 2017.
- Rajendra Bhatia. *Positive definite matrices*. Princeton university press, 2009.
- Silvere Bonnabel. Stochastic gradient descent on Riemannian manifolds. *IEEE Transactions on Automatic Control*, 58(9):2217–2229, 2013.
- Henri Bonnel, Léonard Todjihoundé, and Constantin Udriște. Semivectorial bilevel optimization on Riemannian manifolds. *Journal of Optimization Theory and Applications*, 167(2):464–486, 2015.
- Nicolas Boumal. *An introduction to optimization on smooth manifolds*. Cambridge University Press, 2023.
- Nicolas Boumal and Pierre Absil. RTRMC: A Riemannian trust-region method for low-rank matrix completion. In *Advances in neural information processing systems*, pages 406–414, 2011.
- Nicolas Boumal, Pierre-Antoine Absil, and Coralia Cartis. Global rates of convergence for nonconvex optimization on manifolds. *IMA Journal of Numerical Analysis*, 39(1):1–33, 2018.
- Jerome Bracken and James T McGill. Mathematical programs with optimization problems in the constraints. *Operations Research*, 21(1):37–44, 1973.
- Yang Cai, Michael I Jordan, Tianyi Lin, Argyris Oikonomou, and Emmanouil-Vasileios Vlatakis-Gkaragkounis. Curvature-independent last-iterate convergence for games on Riemannian manifolds. *arXiv preprint arXiv:2306.16617*, 2023.
- Lesi Chen, Jing Xu, and Jingzhao Zhang. On bilevel optimization without lower-level strong convexity. *arXiv preprint arXiv:2301.00712*, 2023a.
- Robert Chen, Brendan Lucier, Yaron Singer, and Vasilis Syrgkanis. Robust optimization for non-convex objectives. *arXiv preprint arXiv:1707.01047*, 2017.
- Tianyi Chen, Yuejiao Sun, and Wotao Yin. Closing the gap: Tighter analysis of alternating stochastic gradient methods for bilevel problems. *Advances in Neural Information Processing Systems*, 34:25294–25307, 2021.
- Tianyi Chen, Yuejiao Sun, Quan Xiao, and Wotao Yin. A single-timescale method for stochastic bilevel optimization. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera, editors, *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pages 2466–2488. PMLR, 28–30 Mar 2022. URL <https://proceedings.mlr.press/v151/chen22e.html>.

- Xuxing Chen, Minhui Huang, Shiqian Ma, and Krishna Balasubramanian. Decentralized stochastic bilevel optimization with improved per-iteration complexity. In *International Conference on Machine Learning*, pages 4641–4671. PMLR, 2023b.
- Xuxing Chen, Minhui Huang, and Shiqian Ma. Decentralized bilevel optimization. *Optimization Letters*, pages 1–65, 2024.
- Anoop Cherian and Suvrit Sra. Riemannian dictionary learning and sparse coding for positive definite matrices. *IEEE transactions on neural networks and learning systems*, 28(12):2859–2871, 2016.
- Constantinos Daskalakis and Ioannis Panageas. The limit points of (optimistic) gradient descent in min-max optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- Justin Domke. Generic methods for optimization-based modeling. In *Artificial Intelligence and Statistics*, pages 318–326. PMLR, 2012.
- Youran Dong, Shiqian Ma, Junfeng Yang, and Chao Yin. A single-loop algorithm for decentralized bilevel optimization. *arXiv preprint arXiv:2311.08945*, 2023.
- Orizon P Ferreira, Mauricio S Louzeiro, and LF4018420 Prudente. Gradient method for optimization on Riemannian manifolds with lower bounded curvature. *SIAM Journal on Optimization*, 29(4):2517–2541, 2019.
- Rémi Flamary, Alain Rakotomamonjy, and Gilles Gasso. Learning constrained task similarities in graph regularized multi-task learning. *Regularization, Optimization, Kernels, and Support Vector Machines*, 103, 2014.
- Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazi, and Massimiliano Pontil. Bilevel programming for hyperparameter optimization and meta-learning. In *International Conference on Machine Learning*, pages 1568–1577. PMLR, 2018.
- Saeed Ghadimi and Mengdi Wang. Approximation methods for bilevel programming. *arXiv preprint arXiv:1802.02246*, 2018.
- Stephen Gould, Basura Fernando, Anoop Cherian, Peter Anderson, Rodrigo Santa Cruz, and Edison Guo. On differentiating parameterized argmin and argmax problems with application to bi-level optimization. *arXiv preprint arXiv:1607.05447*, 2016.
- Riccardo Grazi, Luca Franceschi, Massimiliano Pontil, and Saverio Salzo. On the iteration complexity of hypergradient computation. In *International Conference on Machine Learning*, pages 3748–3758. PMLR, 2020.
- Andi Han, Bamdev Mishra, Pratik Jawanpuria, Pawan Kumar, and Junbin Gao. Riemannian Hamiltonian methods for min-max optimization on manifolds. *SIAM Journal on Optimization*, 33(3):1797–1827, 2023.
- Andi Han, Bamdev Mishra, Pratik Jawanpuria, and Akiko Takeda. A Framework for Bilevel Optimization on Riemannian Manifolds. *arXiv preprint arXiv:2402.03883*, 2024.

- Mehrtash Harandi, Mathieu Salzmann, and Richard Hartley. Dimensionality reduction on SPD manifolds: The emergence of geometry-aware methods. *IEEE transactions on pattern analysis and machine intelligence*, 40(1):48–62, 2017.
- Ken Hayami. Convergence of the conjugate gradient method on singular systems. *arXiv preprint arXiv:1809.00793*, 2018.
- Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang. A two-timescale stochastic algorithm framework for bilevel optimization: Complexity analysis and application to actor-critic. *SIAM Journal on Optimization*, 33(1):147–180, 2023.
- Feihu Huang and Shangqian Gao. Gradient descent ascent for minimax problems on riemannian manifolds. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(7):8466–8476, 2023.
- Kaiyi Ji and Yingbin Liang. Lower bounds and accelerated algorithms for bilevel optimization. *Journal of machine learning research*, 24(22):1–56, 2023.
- Kaiyi Ji, Jason D Lee, Yingbin Liang, and H Vincent Poor. Convergence of meta-learning with task-specific adaptation over partial parameters. *Advances in Neural Information Processing Systems*, 33:11490–11500, 2020.
- Kaiyi Ji, Junjie Yang, and Yingbin Liang. Bilevel optimization: Convergence analysis and enhanced design. In *International Conference on Machine Learning*, pages 4882–4892. PMLR, 2021.
- Chi Jin, Praneeth Netrapalli, and Michael Jordan. What is local optimality in nonconvex-nonconcave minimax optimization? In *International Conference on Machine Learning*, pages 4880–4889. PMLR, 2020.
- Hiroyuki Kasai, Hiroyuki Sato, and Bamdev Mishra. Riemannian stochastic recursive gradient algorithm. In *International Conference on Machine Learning*, pages 2516–2524, 2018.
- Prashant Khanduri, Siliang Zeng, Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang. A near-optimal algorithm for stochastic bilevel optimization via double-momentum. *Advances in neural information processing systems*, 34:30271–30283, 2021.
- Vijay Konda and John Tsitsiklis. Actor-critic algorithms. *Advances in neural information processing systems*, 12, 1999.
- Gautam Kunapuli, Kristin P Bennett, Jing Hu, and Jong-Shi Pang. Classification model selection via bilevel programming. *Optimization Methods & Software*, 23(4):475–489, 2008.
- John M Lee. *Riemannian manifolds: an introduction to curvature*, volume 176. Springer Science & Business Media, 2006a.
- John M Lee. Introduction to smooth manifolds, 2006b.

- Jun Li, Fuxin Li, and Sinisa Todorovic. Efficient riemannian optimization on the stiefel manifold via the cayley transform. In *International Conference on Learning Representations*, 2020a. URL <https://openreview.net/forum?id=HJxV-ANKDH>.
- Junyi Li, Bin Gu, and Heng Huang. Improved bilevel model: Fast and optimal algorithm with theoretical guarantee. *arXiv preprint arXiv:2009.00690*, 2020b.
- Renjie Liao, Yuwen Xiong, Ethan Fetaya, Lisa Zhang, KiJung Yoon, Xaq Pitkow, Raquel Urtasun, and Richard Zemel. Reviving and improving recurrent back-propagation. In *International Conference on Machine Learning*, pages 3082–3091. PMLR, 2018.
- Lizhen Lin, Brian St. Thomas, Hongtu Zhu, and David B Dunson. Extrinsic local regression on manifold-valued data. *Journal of the American Statistical Association*, 112(519): 1261–1273, 2017.
- Lizhen Lin, Drew Lazar, Bayan Sarpabayeva, and David B. Dunson. Robust optimization and inference on manifolds. *arXiv preprint arXiv:2006.06843*, 2020a.
- Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave minimax problems. In *International Conference on Machine Learning*, pages 6083–6093. PMLR, 2020b.
- Risheng Liu, Pan Mu, Xiaoming Yuan, Shangzhi Zeng, and Jin Zhang. A generic first-order algorithmic framework for bi-level programming beyond lower-level singleton. In *International Conference on Machine Learning*, pages 6305–6315. PMLR, 2020.
- Jonathan Lorraine, Paul Vicol, and David Duvenaud. Optimizing millions of hyperparameters by implicit differentiation. In *International Conference on Artificial Intelligence and Statistics*, pages 1540–1552. PMLR, 2020.
- Dougal Maclaurin, David Duvenaud, and Ryan Adams. Gradient-based hyperparameter optimization through reversible learning. In *International conference on machine learning*, pages 2113–2122. PMLR, 2015.
- Bamdev Mishra, Hiroyuki Kasai, Pratik Jawanpuria, and Atul Saroop. A Riemannian gossip approach to subspace learning on Grassmann manifold. *Machine Learning*, pages 1–21, 2019.
- Aryan Mokhtari, Asuman Ozdaglar, and Sarath Pattathil. A unified analysis of extra-gradient and optimistic gradient methods for saddle point problems: Proximal point approach. In *International Conference on Artificial Intelligence and Statistics*, pages 1497–1507. PMLR, 2020.
- Gregory M Moore. *Bilevel programming algorithms for machine learning model selection*. Rensselaer Polytechnic Institute, 2010.
- Yurii Nesterov et al. *Lectures on convex optimization*, volume 137. Springer, 2018.
- Maher Nouiehed, Maziar Sanjabi, Tianjian Huang, Jason D Lee, and Meisam Razaviyayn. Solving a class of non-convex min-max games using iterative first order methods. *Advances in Neural Information Processing Systems*, 32, 2019.

- Takayuki Okuno, Akiko Takeda, Akihiro Kawana, and Motokazu Watanabe. On lp-hyperparameter learning via bilevel nonsmooth optimization. *Journal of Machine Learning Research*, 22(245):1–47, 2021.
- Fabian Pedregosa. Hyperparameter optimization with approximate gradient. In *International conference on machine learning*, pages 737–746. PMLR, 2016.
- Kaare Brandt Petersen, Michael Syskind Pedersen, et al. The matrix cookbook. *Technical University of Denmark*, 7(15):510, 2008.
- Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. Meta-learning with implicit gradients. *Advances in neural information processing systems*, 32, 2019.
- Amirreza Shaban, Ching-An Cheng, Nathan Hatch, and Byron Boots. Truncated back-propagation for bilevel optimization. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1723–1732. PMLR, 2019.
- Chenggen Shi, Jie Lu, and Guangquan Zhang. An extended Kuhn–Tucker approach for linear bilevel programming. *Applied Mathematics and Computation*, 162(1):51–63, 2005.
- Suvrit Sra and Reshad Hosseini. Conic geometric optimization on the manifold of positive definite matrices. *SIAM Journal on Optimization*, 25(1):713–739, 2015.
- Heinrich von Stackelberg and Alan T Peacock. The theory of the market economy. (*No Title*), 1952.
- Ju Sun, Qing Qu, and John Wright. Complete dictionary recovery over the sphere ii: Recovery by Riemannian trust-region method. *IEEE Transactions on Information Theory*, 63(2):885–914, 2016.
- Ju Sun, Qing Qu, and John Wright. A geometric analysis of phase retrieval. *Foundations of Computational Mathematics*, 18(5):1131–1198, 2018.
- Davoud Ataee Tarzanagh, Mingchen Li, Christos Thrampoulidis, and Samet Oymak. Fednest: Federated bilevel, minimax, and compositional optimization. In *International Conference on Machine Learning*, pages 21146–21179. PMLR, 2022.
- Nilesh Tripuraneni, Nicolas Flammarion, Francis Bach, and Michael I Jordan. Averaging stochastic gradient descent on Riemannian manifolds. In *Conference On Learning Theory*, pages 650–687, 2018.
- Loring W Tu. *An Introduction to Manifolds*. Springer Science & Universitext, 2011.
- Bart Vandereycken. Low-rank matrix completion by Riemannian optimization. *SIAM Journal on Optimization*, 23(2):1214–1236, 2013.
- Melanie Weber and Suvrit Sra. Projection-free nonconvex stochastic optimization on Riemannian manifolds. *IMA Journal of Numerical Analysis*, 42(4):3241–3271, 2022.
- Melanie Weber and Suvrit Sra. Riemannian optimization via Frank-Wolfe methods. *Mathematical Programming*, 199(1-2):525–556, 2023.

- Yifan Yang, Peiyao Xiao, and Kaiyi Ji. Achieving $\mathcal{O}(\epsilon^{-1.5})$ complexity in Hessian/Jacobian-free stochastic bilevel optimization. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=0zjBohmLvE>.
- TaeHo Yoon and Ernest K Ryu. Accelerated algorithms for smooth convex-concave minimax problems with $\mathcal{O}(1/k^2)$ rate on squared gradient norm. In *International Conference on Machine Learning*, pages 12098–12109. PMLR, 2021.
- Tong Yu and Hong Zhu. Hyper-parameter optimization: A review of algorithms and applications. *arXiv preprint arXiv:2003.05689*, 2020.
- Hongyi Zhang and Suvrit Sra. First-order methods for geodesically convex optimization. In *Conference on Learning Theory*, pages 1617–1638. PMLR, 2016.
- Hongyi Zhang, Sashank J Reddi, and Suvrit Sra. Riemannian SVRG: Fast stochastic optimization on Riemannian manifolds. *Advances in Neural Information Processing Systems*, 29:4592–4600, 2016.
- Pan Zhou, Xiaotong Yuan, Shuicheng Yan, and Jiashi Feng. Faster first-order methods for stochastic non-convex optimization on Riemannian manifolds. *IEEE transactions on pattern analysis and machine intelligence*, 2019.