

# Stochastic Zeroth-order Riemannian Derivative Estimation and Optimization

Jiaxiang Li\*

Krishnakumar Balasubramanian<sup>†</sup>

Shiqian Ma<sup>‡</sup>

January 6, 2021

## Abstract

We consider stochastic zeroth-order optimization over Riemannian submanifolds embedded in Euclidean space, where the task is to solve Riemannian optimization problem with only noisy objective function evaluations. Towards this, our main contribution is to propose estimators of the Riemannian gradient and Hessian from noisy objective function evaluations, based on a Riemannian version of the Gaussian smoothing technique. The proposed estimators overcome the difficulty of the non-linearity of the manifold constraint and the issues that arise in using Euclidean Gaussian smoothing techniques when the function is defined only over the manifold. We use the proposed estimators to solve Riemannian optimization problems in the following settings for the objective function: (i) stochastic and gradient-Lipschitz (in both nonconvex and geodesic convex settings), (ii) sum of gradient-Lipschitz and non-smooth functions, and (iii) Hessian-Lipschitz. For these settings, we analyze the oracle complexity of our algorithms to obtain appropriately defined notions of  $\epsilon$ -stationary point or  $\epsilon$ -approximate local minimizer. Notably, our complexities are independent of the dimension of the ambient Euclidean space and depend only on the intrinsic dimension of the manifold under consideration. We demonstrate the applicability of our algorithms by simulation results and real-world applications on black-box stiffness control for robotics and black-box attacks to neural networks.

## 1 Introduction

Consider the following Riemannian optimization problem:

$$\min f(x) + h(x), \text{ s.t., } x \in \mathcal{M}, \quad (1.1)$$

where  $\mathcal{M}$  is a Riemannian submanifold embedded in  $\mathbb{R}^n$ ,  $f : \mathcal{M} \rightarrow \mathbb{R}$  is a smooth and possibly nonconvex function, and  $h : \mathbb{R}^n \rightarrow \mathbb{R}$  is a convex and nonsmooth function. Here, convexity and smoothness are interpreted as the function is being considered in the ambient Euclidean space. Iterative algorithms for solving (1.1) usually require the gradient and Hessian information of the objective function. However, in many applications, the analytical form of the function  $f$  (or  $h$ ) and its gradient are not available, and we can only obtain noisy function evaluations via a zeroth-order oracle. This setting, termed as *stochastic zeroth-order Riemannian optimization*, generalizes stochastic zeroth-order Euclidean optimization (i.e., when  $\mathcal{M} \equiv \mathbb{R}^n$  in (1.1)), a topic which goes back to the early works of [Mat65, NM65, NY83] in the 1960's; see also [CSV09, AH17, LMW19] for recent books and surveys.

---

\*Department of Mathematics, University of California, Davis. [jxjli@ucdavis.edu](mailto:jxjli@ucdavis.edu).

<sup>†</sup>Department of Statistics, University of California, Davis. [kbala@ucdavis.edu](mailto:kbala@ucdavis.edu).

<sup>‡</sup>Department of Mathematics, University of California, Davis. [sqma@ucdavis.edu](mailto:sqma@ucdavis.edu).

In the Euclidean setting, two popular techniques for estimating the gradient from (noisy) function queries include the finite-differences method [Spa05] and the Gaussian smoothing techniques [NY83]. Earlier works in this setting focused on using the estimated gradient to obtain asymptotic convergence rates of iterative optimization algorithms. Recently, obtaining non-asymptotic guarantees on the oracle complexity of stochastic zeroth-order optimization has been of great interest. Towards that, [Nes11, NS17] analyzed the Gaussian smoothing technique for estimating the Euclidean gradient from noisy function evaluations and proved that for unconstrained convex minimization, one needs  $O(n^2/\epsilon^2)$  noisy function evaluations to obtain an  $\epsilon$ -optimal solution. This complexity was improved by [GL13] to  $O(n/\epsilon^2)$  when the objective function is further assumed to be gradient-smooth. Note that this oracle complexity depends linearly on the problem dimension  $n$  and it was proved that the linear dependency on  $n$  is unavoidable [JNR12, DJWW15]. Nonconvex and smooth setting was also considered in [GL13]. In particular, now assuming  $h \equiv 0$  and  $\mathcal{M} \equiv \mathbb{R}^n$  in (1.1), it was shown that the number of function evaluations for obtaining an  $\epsilon$ -stationary point  $\bar{x}$  (i.e.,  $\mathbb{E}\|\nabla f(\bar{x})\| \leq \epsilon$ ), is  $O(n/\epsilon^4)$ .

In the Riemannian manifold setting, however, a main challenge in designing and analyzing zeroth-order algorithms is the lack of availability of theoretically sound methods to estimate the Riemannian gradients and Hessians from (noisy) function evaluations. To this end, our main contribution in this work is to construct estimators of the Riemannian gradient and Hessian from noisy function evaluations, based on a modified Gaussian smoothing technique from [NS17] and [BG19]. The main difficulty addressed here is that the gradient and Hessian estimator in [NS17] and [BG19] respectively, require computing  $f(x + \nu u)$ , for some parameter  $\nu > 0$  and an  $n$ -dimensional standard Gaussian vector  $u$ . However, the point  $x + \nu u$  may not necessarily lie on the manifold  $\mathcal{M}$ . To resolve this issue, we propose an estimator based on the smoothing technique and sampling Gaussian random vectors on the tangent space of the manifold  $\mathcal{M}$ . We then use the developed methodology to design stochastic zeroth-order algorithms for solving (1.1) with oracle complexities that depend only on the manifold dimension  $d$ , and independent of the ambient Euclidean dimension  $n$ .

## 1.1 Related Works

As mentioned previously, zeroth-order optimization has a long history; we refer the reader to [CSV09, AH17, LMW19] for more details. The oracle complexity of methods from the above works are at least linear in terms of their dependence on dimensionality. Recent works in this field have been focusing on stochastic zeroth-order optimization in high-dimensions [WDBS18, GKK<sup>+</sup>19, BG19, CMYZ20]. Assuming a sparse structure (for example, the function being optimized depends only on  $s$  of the  $n$  coordinates), the above works have shown that the oracle complexity of zeroth-order optimization depends only poly-logarithmically on the dimension  $n$ , and it has a linear dependency only on the sparsity parameter  $s$ , which is typically small compared to  $n$  in several applications. Compared to these works, we assume a manifold structure on the function being optimized and obtain oracle complexities that depend only on the manifold dimension and independent of the ambient Euclidean dimension.

Apart from the above, *Bayesian optimization* is yet another popular class of methods for optimizing functions based on noisy function values. This approach aims at finding the global minimizer by enforcing a Gaussian process prior on the space of function being optimized and using Bayesian sampling techniques. We refer the reader to [Moc94, Moc12, SSW<sup>+</sup>15, Fra18] for an overview of such techniques in the Euclidean settings and their applications to a variety of fields including robotics, recommender systems, preference learning and hyperparameter tuning. A common limitation of the above algorithms is that they are usually not scalable well to solve high-dimensional problems. Recent developments on Bayesian optimization for high-dimensional

problems include [LKPS16, WHZ<sup>+</sup>16, MK18, RSBC18, WFT20] where people considered zeroth-order optimization with structured functions (for example, sparse or additive functions), and developed Bayesian optimization algorithms and related analysis. Very recently, [OGW18, JRCB20, JR20] considered heuristic Bayesian optimization algorithms for function defined over non-Euclidean domains, including Riemannian domains, without any theoretical analysis.

Riemannian optimization in the first or second-order setting has drawn a lot of attention recently due to its applications in various fields, including low-rank matrix completion [BA11, Van13], phase retrieval [BEB17, SQW18], dictionary learning [CS16, SQW16], dimensionality reduction [HSH17, TFBJ18, MKJS19] and manifold regression [LSTZD17, LLS20]. For smooth Riemannian optimization, i.e.,  $h \equiv 0$  in (1.1), it was shown that Riemannian gradient descent method require  $\mathcal{O}(1/\epsilon^2)$  iterations to converge to an  $\epsilon$ -stationary point [BAC18]. Stochastic algorithms were also studied for smooth Riemannian optimization [Bon13, ZYYF19, WS19, ZRS16, KSM18, ZYYF19, WS19]. In particular, using the SPIDER variance reduction technique, [ZYYF19] proved that  $\mathcal{O}(1/\epsilon^3)$  oracle calls are required to obtain an  $\epsilon$ -stationary point in expectation. When the function  $f$  takes a finite-sum structure, the Riemannian SVRG [ZRS16] achieves  $\epsilon$ -stationary solution with  $\mathcal{O}(k^{2/3}/\epsilon^2)$  oracle calls where  $k$  is number of summands. When the nonsmooth function  $h$  presents in (1.1), Riemannian sub-gradient methods (RSGM) are widely used [BSBA14, LCD<sup>+</sup>19] and they require  $\mathcal{O}(1/\epsilon^4)$  iterations. ADMM for solving (1.1) has also been studied [KGB16, LO14], but they usually lack convergence guarantee, while the analysis presented in [ZMZ20] requires some strong assumptions. The recently proposed manifold proximal gradient method (ManPG) [CMMCSZ20] for solving (1.1) requires  $\mathcal{O}(1/\epsilon^2)$  number of iterations to find an  $\epsilon$ -stationary solution. Variants of ManPG such as ManPPA [CDMS20], ManPL [WLC<sup>+</sup>20] and stochastic ManPG [WMX20] have also been studied. Note that none of these works considers the zeroth-order setting. Recently, there are some attempts on stochastic zeroth-order Riemannian optimization [CSA15, FT19], but they are mostly heuristics and do not have any rigorous convergence guarantees.

## 1.2 Motivating Applications

Our motivation for developing a theoretical framework for stochastic zeroth-order Riemannian optimization is due to several important emerging applications; see, e.g., [CSA15, MHST17, YCL19, JRCB20, Kac20]. Below, we discuss two concrete examples, which we will revisit in Section 4.2, to illustrate the applicability of the methods developed in this work. We also briefly discuss a third application in topological data analysis, and numerical experiments on this application will be conducted in a future work, as it is more involved and beyond the scope of this paper.

### 1.2.1 Black-box Stiffness Control for Robotics

Our first motivating application is from the field of robotics. It has become increasingly common to use zeroth-order optimization techniques to optimize control parameter and policies in robotics [MHB<sup>+</sup>16, DET17, YCL19]. This is because that the cost functions being optimized in robotics are not available in a closed form as a function of the control parameter. Invariably for a given choice of control parameter, the cost function needs to be evaluated through a real-world experiment on a given robot or through simulation. Recently, domain knowledge has been used as constraints on the control parameter space, among which a common choice is the geometry-aware constraint. For example, control parameters like stiffness, inertia and manipulability lie on the positive semidefinite manifold, orthogonal group and unit sphere, respectively. Hence, there is a need to develop zeroth-order optimization methods over the manifolds to optimize the above mentioned control parameters [JRCB20].

### 1.2.2 Zeroth-order Attacks on Deep Neural Networks (DNNs)

Our second motivating application is based on developing black-box attacks to DNNs. Despite the recent success of DNNs, studies have shown that they are vulnerable to adversarial attacks: even a well-trained DNN could completely misclassify a slightly perturbed version of the original image (which is undetectable to the human eyes); see, e.g., [SZS<sup>+</sup>13, GSS14]. As a result, it is extremely important on the one hand to come up with methods to train DNNs that are robust to adversarial attacks, and on the other hand to develop efficient attacks on DNNs with the goal being to make them misclassify. In practice, as the architecture of the DNN is not known to the attacker, several works, for example, [CZS<sup>+</sup>17, TTC<sup>+</sup>19, CLC<sup>+</sup>18], use zeroth-order optimization algorithms for designing adversarial attacks. However, existing works have an inherent drawback—the perturbed training example designed to fool the DNN is not in the same domain as the original training data. For example, despite the fact that natural images typically lie on a manifold [WS04], the perturbations are not constrained to lie on the same manifold. This naturally motivates us to use zeroth-order Riemannian optimization methods to design adversarial examples to fool DNNs, which at the same time preserves the manifold structures in the training data.

### 1.2.3 Black-box Methods for Topological Dimension Reduction

The third motivating example is from the field of dimension reduction, a popular class of techniques for reducing the dimension of high-dimensional unstructured data for feature extraction and visualization. There exists a variety of methods for this task; we refer the interested reader to [LV07, Bur10] for more details. However, a majority of the existing techniques are based on *geometric* motivations. Recently, there has been a growing literature on using *topological* information for performing data analysis [CM17, MHM18, RB19]. One such method is a dimension reduction technique called Persistent Homology-Based Projection Pursuit [Kac20]. Roughly speaking, given a point-cloud data set with cardinality  $n$  and dimension  $m$  (i.e., a matrix  $X \in \mathbb{R}^{m \times n}$ ), persistence homology refers to developing a multi-scale characterization of topologically invariant features available in the data. Such information is summarized in terms of the so-called persistence diagram,  $D(X)$ , which is a multiset of points in a two-dimensional plane. The idea in [Kac20] is to obtain a transformation  $P^\top \in \mathbb{R}^{p \times m}$ , with  $p \ll m$ , such that the topological summaries of the original dataset  $X$  and the reduced dimensional dataset  $P^\top X$  are close to each other; that is, the persistence diagram  $D(X)$  and  $D(P^\top X)$  are close in the 2-Wasserstein distance. The problem is then formulated as (informally speaking),

$$\min_{\{P \in \mathbb{R}^{m \times p} : P^\top P = I\}} W_2(D(X), D(P^\top X)),$$

which is an optimization problem over the Stiefel manifold (see Section 2.1 for more details). It turns out that calculating the gradient of the above objective function is highly non-trivial and computationally expensive [LOT19]. However, evaluating the objective function for various value of the matrix  $P$  is relatively less expensive. Hence, this serves as yet another problem in which the methodology developed in this work could be applied naturally.

## 1.3 Main Contributions

We now summarize our main contributions.

1. In Section 2, we propose the (stochastic) zeroth-order Riemannian gradient (2.5) and Hessian (2.14) estimators, which addresses the infeasibility issue of the sampling for the case of derivative-free optimization over manifolds.

ALGORITHM	STRUCTURE	ITERATION COMPLEXITY	ORACLE COMPLEXITY
ZO-RGD	SMOOTH	$\mathcal{O}(d/\epsilon^2)$	$\mathcal{O}(d/\epsilon^2)$
ZO-RSGD	SMOOTH, STOCHASTIC	$\mathcal{O}(1/\epsilon^2)$	$\mathcal{O}(d/\epsilon^4)$
ZO-RSGD	SMOOTH, STOCHASTIC, GEO-CONVEX	$\mathcal{O}(1/\epsilon)$	$\mathcal{O}(d/\epsilon^2)$
ZO-SManPG	NONSMOOTH STOCHASTIC	$\mathcal{O}(1/\epsilon^2)$	$\mathcal{O}(d/\epsilon^4)$
ZO-RSCRN	LIPSCHITZ HESSIAN STOCHASTIC	$\mathcal{O}(1/\epsilon^{1.5})$	$\mathcal{O}(d/\epsilon^{3.5} + d^4/\epsilon^{2.5})$

Table 1: Summary of the convergence results proved in this paper. For all but the ZO-RSCRN algorithm, the reported complexities correspond to  $\epsilon$ -stationary solution; for the ZO-RSCRN algorithm the complexities correspond to  $\epsilon$ -local minimizers. Here,  $d$  is the intrinsic dimension of the manifold  $\mathcal{M}$ . Furthermore, Iteration complexity refers to the number of iterations and oracle complexity refers to the number of calls to the (stochastic) zeroth-order oracle.

2. In Section 3, we demonstrate the applicability of the developed estimators for stochastic zeroth-order Riemannian optimization, as listed below. A summary of these results is given in Table 1. To the best of our knowledge, our results are the first complexity results for stochastic zeroth-order Riemannian optimization.
  - When  $h(x) \equiv 0$  and the exact function evaluations of  $f(x)$  are obtainable, we propose a zeroth-order Riemannian gradient descent method (ZO-RGD) and provide its oracle complexity for obtaining an  $\epsilon$ -stationary point of (1.1) (see Theorem 3.1).
  - When  $h(x) \equiv 0$  and  $f(x) = \mathbb{E}_\xi[F(x, \xi)]$ , we propose a zeroth-order Riemannian stochastic gradient descent method (ZO-RSGD). We analyze its oracle complexity under two different settings (see Theorems 3.2 and A.1).
  - When  $h(x)$  is convex and nonsmooth, we propose a zeroth-order stochastic Riemannian proximal gradient method (ZO-SManPG) and provide its oracle complexity for obtaining an  $\epsilon$ -stationary point of (1.1) (see Theorem 3.3).
  - When  $h(x) \equiv 0$  and  $f(x) = \mathbb{E}_\xi[F(x, \xi)]$ , where  $F(x, \xi)$  satisfies certain Lipschitz Riemannian Hessian property, we propose a zeroth-order Riemannian stochastic cubic regularized Newton’s method (ZO-RSCRN) that provably converges to an  $\epsilon$ -approximate local minimizer (see Theorem 3.4).
3. In Section 4, we provide experimental results on simulated data to quantify the performance of our methods. We then demonstrate the applicability of our methods to the problem of black-box attacks to deep neural networks and robotics.

## 2 Preliminaries and Methodology

We start this section with a brief review of basics of Riemannian optimization. We then introduce our stochastic zeroth-order Riemannian gradient and Hessian estimators, and provide bias and moment bounds quantifying the accuracy of the proposed estimators, which will be useful for the convergence analysis later.

## 2.1 Basics of Riemannian Optimization

Let  $\mathcal{M} \subset \mathbb{R}^n$  be a differentiable embedded submanifold. We have the following definition for the tangent space.

**Definition 2.1** (Tangent space). *Consider a manifold  $\mathcal{M}$  embedded in a Euclidean space. For any  $x \in \mathcal{M}$ , the tangent space  $T_x\mathcal{M}$  at  $x$  is a linear subspace that consists of the derivatives of all differentiable curves on  $\mathcal{M}$  passing through  $x$ :*

$$T_x\mathcal{M} = \{\gamma'(0) : \gamma(0) = x, \gamma([- \delta, \delta]) \subset \mathcal{M} \text{ for some } \delta > 0, \gamma \text{ is differentiable}\}. \quad (2.1)$$

The manifold  $\mathcal{M}$  is a Riemannian manifold if it is equipped with an inner product on the tangent space,  $\langle \cdot, \cdot \rangle_x : T_x\mathcal{M} \times T_x\mathcal{M} \rightarrow \mathbb{R}$ , that varies smoothly on  $\mathcal{M}$ . We also introduce the concept of the dimension of a manifold.

**Definition 2.2** (Dimension of a manifold [AMS09]). *The dimension of the manifold  $\mathcal{M}$ , denoted as  $d$ , is the dimension of the Euclidean space that the manifold is locally homeomorphic to. In particular, the dimension of the tangent space is always equal to the dimension of the manifold.*

As an example, consider the Stiefel manifold  $\mathcal{M} = \text{St}(n, p) := \{X \in \mathbb{R}^{n \times p} : X^\top X = I_p\}$ . The tangent space of  $\text{St}(n, p)$  is given by  $T_X\mathcal{M} = \{Y \in \mathbb{R}^{n \times p} : X^\top Y + Y^\top X = 0\}$ . Hence, the dimension of the Stiefel manifold is  $np - \frac{1}{2}p(p+1)$ . Note that the dimension of the manifold could be significantly less than the dimension of the ambient Euclidean space. Yet another example is the manifold of low-rank matrices [Van13]. We now introduce the concept of a Riemannian gradient.

**Definition 2.3** (Riemannian Gradient). *Suppose  $f$  is a smooth function on  $\mathcal{M}$ . The Riemannian gradient  $\text{grad}f(x)$  is a vector in  $T_x\mathcal{M}$  satisfying  $\left. \frac{d(f(\gamma(t)))}{dt} \right|_{t=0} = \langle v, \text{grad}f(x) \rangle_x$  for any  $v \in T_x\mathcal{M}$ , where  $\gamma(t)$  is a curve as described in (2.1).*

Recall that in the Euclidean setting, a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $L$ -smooth, if it satisfies  $|f(y) - f(x) - \langle \nabla f(x), y - x \rangle| \leq \frac{L}{2} \|x - y\|^2$ , for all  $x, y \in \mathbb{R}^n$ . We now present the Riemannian counterpart of  $L$ -smooth functions. To do so, first we need the definition of retraction for a given  $x \in \mathcal{M}$ .

**Definition 2.4** (Retraction). *A retraction mapping  $R_x$  is a smooth mapping from  $T_x\mathcal{M}$  to  $\mathcal{M}$  such that:  $R_x(0) = x$ , where  $0$  is the zero element of  $T_x\mathcal{M}$ , and the differential of  $R_x$  at  $0$  is an identity mapping, i.e.,  $\left. \frac{dR_x(t\eta)}{dt} \right|_{t=0} = \eta, \forall \eta \in T_x\mathcal{M}$ . In particular, the exponential mapping  $\text{Exp}_x$  is a retraction that generates geodesics.*

**Assumption 2.1** ( $L$ -retraction-smoothness). *There exists  $L_g \geq 0$  such that the following inequality holds for function  $f$  in (1.1):*

$$|f(R_x(\eta)) - f(x) - \langle \text{grad}f(x), \eta \rangle_x| \leq \frac{L_g}{2} \|\eta\|^2, \forall x \in \mathcal{M}, \eta \in T_x\mathcal{M}. \quad (2.2)$$

Assumption 2.1 is also known as the restricted Lipschitz-type gradient for pullback function  $\hat{f}_x(\eta) := f(R_x(\eta))$  [BAC18]. The condition required in [BAC18] is weaker because it only requires Eq. (2.2) to hold for  $\|\eta\|_x \leq \rho_x$ , where constant  $\rho_x > 0$ . In our convergence analysis, we need this assumption to be held for all  $\eta \in T_x\mathcal{M}$ , i.e.,  $\rho_x = \infty$ . This assumption is satisfied when the manifold  $\mathcal{M}$  is a compact submanifold of  $\mathbb{R}^n$ , the retraction  $R_x$  is globally defined<sup>1</sup> and function  $f$  is  $L$ -smooth

<sup>1</sup>If the manifold is compact, then the exponential mapping  $\text{Exp}_x$  is already globally defined. This is known as the Hopf-Rinow theorem [Car92].

in the Euclidean sense; we refer the reader to [BAC18] for more details. We also emphasize that Assumption 2.1 is weaker than the geodesic smoothness assumption defined in [ZS16]. The geodesic smoothness states that,  $\forall \eta \in \mathcal{M}$ ,  $f(\text{Exp}_x(\eta)) \leq f(x) + \langle g_x, \eta \rangle_x + L_g d^2(x, \text{Exp}_x(\eta))/2$ , where  $g_x$  is a subgradient of  $f$ ,  $d(\cdot, \cdot)$  represents the geodesic distance. Such a condition is stronger than our Assumption 2.1, in the sense that, if the retraction is the exponential mapping, then geodesic smoothness implies the  $L$ -retraction-smoothness with the same parameter  $L_g$  [BFM17].

Throughout this paper, we consider the Riemannian metric on  $\mathcal{M}$  that is induced from the Euclidean inner product; i.e.  $\langle \cdot, \cdot \rangle_x = \langle \cdot, \cdot \rangle$ ,  $\forall x \in \mathcal{M}^2$ . Using this Riemannian metric, the Riemannian gradient of a function is simply the projection of its Euclidean gradient onto the tangent space:

$$\text{grad}f(x) = \text{Proj}_{T_x\mathcal{M}}(\nabla f(x)). \quad (2.3)$$

We also present the definition of Riemannian Hessian for embedded submanifolds, which will be used in Section 3.4 about cubic regularized Newton's method.

**Definition 2.5** (Riemannian Hessian [ZZ18]). *Suppose  $\mathcal{M}$  is an embedded submanifold of  $\mathbb{R}^n$ . The Riemannian Hessian is defined as*

$$\text{Hess}f(x)[\eta] = \text{Proj}_{T_x\mathcal{M}}(D\text{grad}f(x)[\eta]), \forall x \in \mathcal{M}, \eta \in T_x\mathcal{M}, \quad (2.4)$$

where  $D\text{grad}f(x)[\eta]$  is the common differential, i.e.,  $D\text{grad}f(x)[\eta] = (J\text{grad}f(x))[\eta]$ , where  $J$  is the Jacobian of the gradient mapping.

## 2.2 The Zeroth-order Riemannian Gradient Estimator

Recall that in the Euclidean setting, Nesterov and Spokoiny [NS17] analyzed the Gaussian smoothing based zeroth-order gradient estimator. However, as that estimator requires function evaluations outside of the manifold to be well-defined, it is not directly applicable for the Riemannian setting. To address this issue, we introduce our stochastic zeroth-order Riemannian gradient estimator below.

**Definition 2.6** (Zeroth-Order Riemannian Gradient). *Generate  $u = Pu_0 \in T_x\mathcal{M}$ , where  $u_0 \sim \mathcal{N}(0, I_n)$  in  $\mathbb{R}^n$ , and  $P \in \mathbb{R}^{n \times n}$  is the orthogonal projection matrix onto  $T_x\mathcal{M}$ . Therefore  $u$  follows the standard normal distribution  $\mathcal{N}(0, PP^\top)$  on the tangent space in the sense that, all the eigenvalues of the covariance matrix  $PP^\top$  are either 0 (eigenvectors orthogonal to the tangent plane) or 1 (eigenvectors embedded in the tangent plane). The zeroth-order Riemannian gradient estimator is defined as*

$$g_\mu(x) = \frac{f(R_x(\mu u)) - f(x)}{\mu} u = \frac{f(R_x(\mu Pu_0)) - f(x)}{\mu} Pu_0. \quad (2.5)$$

Note that the projection  $P$  is easy to compute for commonly used manifolds. For example, for the Stiefel manifold  $\mathcal{M}$ , the projection is given by  $\text{Proj}_{T_x\mathcal{M}}(Y) = (I - XX^\top)Y + X \text{skew}(X^\top Y)$ , where  $\text{skew}(A) := (A - A^\top)/2$  (see [AMS09]).

**Remark 2.1.** *In this work, we assume that the function  $f$  is defined on submanifolds embedded in Euclidean space, so that it is efficient to sample from the associated tangent space, as discussed above; see also [DHS<sup>+</sup>13]. We remark that the above gradient estimation methodology is more generally applicable to other manifolds. However, the generality comes at the cost of practical applicability as it is not an easy task to efficiently sample Gaussian random objects on the tangent space of general manifolds; see [Hsu02] for more details.*

---

<sup>2</sup>If the manifold is not an embedded submanifold of some Euclidean space, then we cannot have an induced Riemannian metric. In this case, the convergence result is not affected, though it would cause implementation difficulties.

We now discuss some differences between the zeroth-order gradient estimators in the Euclidean setting [NS17] and the Riemannian setting (2.5). In the Euclidean case, the zeroth-order gradient estimator can be viewed as estimating the gradient of the Gaussian smoothed function,  $f_\mu(x) = \frac{1}{\kappa} \int_{\mathbb{R}^n} f(x + \mu u) e^{-\frac{1}{2}\|u\|^2} du$ , because  $\nabla f_\mu(x) = \mathbb{E}_u(g_\mu(x)) = \frac{1}{\kappa} \int_{\mathbb{R}^n} \frac{f(x+\mu u) - f(x)}{\mu} u e^{-\frac{1}{2}\|u\|^2} du$ , where  $\kappa$  is the normalization constant for Gaussian. This was also observed as an instantiation of Gaussian Stein's identity [BG19]. However, this observation is no longer true in the Riemannian setting, as we incorporate the retraction operator when evaluating  $g_\mu$ , and this forces us to seek for a direct evaluation of  $\mathbb{E}_u(g_\mu(x))$ , instead of utilizing properties of the smoothed function  $f_\mu$ . We also remark that,  $g_\mu(x)$  is a biased estimator of  $\text{grad}f(x)$ . The difference between them can be bounded as in Proposition 2.1. Some intermediate results for this purpose are as follows.

**Lemma 2.1.** *Suppose  $\mathcal{X}$  is a  $d$ -dimensional subspace of  $\mathbb{R}^n$ , with orthogonal projection matrix  $P \in \mathbb{R}^{n \times n}$ .  $u_0$  follows a standard normal distribution  $\mathcal{N}(0, I_n)$ , and  $u = Pu_0$  is the orthogonal projection of  $u_0$  onto the subspace  $\mathcal{X}$ . Then  $\forall x \in \mathcal{X}$ , we have*

$$x = \frac{1}{\kappa} \int_{\mathbb{R}^n} \langle x, u \rangle u e^{-\frac{1}{2}\|u_0\|^2} du_0, \quad \text{and} \quad \|x\|^2 = \frac{1}{\kappa} \int_{\mathbb{R}^n} \langle x, u \rangle^2 e^{-\frac{1}{2}\|u_0\|^2} du_0, \quad (2.6)$$

where  $\kappa$  is the constant for normal density function:  $\kappa := \int_{\mathbb{R}^n} e^{-\frac{1}{2}\|u\|^2} du = (2\pi)^{n/2}$ .

*Proof.* Proof of Lemma 2.1 By the definition of covariance matrix, we have  $\frac{1}{\kappa} \int_{\mathbb{R}^n} u_0 u_0^\top e^{-\frac{1}{2}\|u_0\|^2} du_0 = I_n$ . Since  $\langle x, u \rangle = \langle x, u_0 \rangle$ ,  $\forall x \in \mathcal{X}$ , we have

$$\frac{1}{\kappa} \int_{\mathbb{R}^n} \langle x, u \rangle u_0 e^{-\frac{1}{2}\|u_0\|^2} du_0 = x, \quad (2.7)$$

which implies  $\frac{1}{\kappa} \int_{\mathbb{R}^n} \langle x, u \rangle u e^{-\frac{1}{2}\|u_0\|^2} du_0 = Px = x$ . Similarly, taking inner product with  $x$  on both sides of Eq. (2.7), we have  $\|x\|^2 = \frac{1}{\kappa} \int_{\mathbb{R}^n} \langle x, u \rangle^2 e^{-\frac{1}{2}\|u_0\|^2} du_0$ .  $\square$

The following bound for the moments of normal distribution is restated without proof.

**Lemma 2.2.** [NS17] *Suppose  $u \sim \mathcal{N}(0, I_n)$  is a standard normal distribution. Then for all integers  $p \geq 2$ , we have  $M_p := \mathbb{E}_u(\|u\|^p) \leq (n+p)^{p/2}$ .*

**Corollary 2.1.** *For  $u_0 \sim \mathcal{N}(0, I_n)$  and  $u = Pu_0$ , where  $P \in \mathbb{R}^{n \times n}$  is the orthogonal projection matrix onto a  $d$  dimensional subspace  $\mathcal{X}$  of  $\mathbb{R}^n$ , we have  $\mathbb{E}_{u_0}(\|u\|^p) \leq (d+p)^{p/2}$ .*

*Proof.* Proof of Corollary 2.1 Assume the eigen-decomposition of  $P$  is  $P = Q^\top \Lambda Q$ , where  $Q$  is an unitary matrix and  $\Lambda$  is a diagonal matrix with the leading  $d$  diagonal entries being 1 and other diagonal entries being 0. Denote  $\tilde{u} = Qu_0 \sim \mathcal{N}(0, I_n)$ , then  $\Lambda \tilde{u} = (\tilde{u}_1, \dots, \tilde{u}_d, 0, \dots, 0)$ . Since  $u = Q^\top \Lambda \tilde{u}$  has the same distribution as  $\Lambda \tilde{u}$ , we have  $\mathbb{E}\|u\|^p = \mathbb{E}\|(\tilde{u}_1, \dots, \tilde{u}_d, 0, \dots, 0)\|^p \leq (d+p)^{p/2}$ , by Lemma 2.2.  $\square$

Now we provide the bounds on the error of our gradient estimator  $g_\mu(x)$  (2.5). Recall that  $d$  denotes the dimension of the manifold  $\mathcal{M}$ .

**Proposition 2.1.** *Under Assumption 2.1, we have*

- (a)  $\|\mathbb{E}_{u_0}(g_\mu(x)) - \text{grad}f(x)\| \leq \frac{\mu L_g}{2} (d+3)^{3/2}$ ,
- (b)  $\|\text{grad}f(x)\|^2 \leq 2\|\mathbb{E}_{u_0}(g_\mu(x))\|^2 + \frac{\mu^2}{2} L_g (d+6)^3$ ,
- (c)  $\mathbb{E}_{u_0}(\|g_\mu(x)\|^2) \leq \frac{\mu^2}{2} L_g^2 (d+6)^3 + 2(d+4)\|\text{grad}f(x)\|^2$ .



*Proof.* Proof of Proposition 2.1 For part (a), since

$$\mathbb{E}(g_\mu(x)) - \text{grad}f(x) = \frac{1}{\kappa} \int_{\mathbb{R}^n} \left( \frac{f(R_x(\mu u)) - f(x)}{\mu} - \langle \text{grad}f(x), u \rangle \right) u e^{-\frac{1}{2}\|u_0\|^2} du_0,$$

we have

$$\begin{aligned} & \|\mathbb{E}(g_\mu(x)) - \text{grad}f(x)\| \\ &= \left\| \frac{1}{\mu\kappa} \int_{\mathbb{R}^n} (f(R_x(\mu u)) - f(x) - \langle \text{grad}f(x), \mu u \rangle) u e^{-\frac{1}{2}\|u_0\|^2} du_0 \right\| \\ &\leq \frac{1}{\mu\kappa} \int_{\mathbb{R}^n} \frac{L_g}{2} \|\mu u\|^2 \|u\| e^{-\frac{1}{2}\|u_0\|^2} du_0 = \frac{\mu L_g}{2\kappa} \int_{\mathbb{R}^n} \|u\|^3 e^{-\frac{1}{2}\|u_0\|^2} du_0 \leq \frac{\mu L_g}{2} (d+3)^{3/2}, \end{aligned}$$

where the first inequality is by due to (2.2), and the last inequality is from Corollary 2.1. This completes the proof of part (a).

To prove part (b), note that

$$\begin{aligned} \|\text{grad}f(x)\|^2 &= \left\| \frac{1}{\kappa} \int_{\mathbb{R}^n} \langle \text{grad}f(x), u \rangle u e^{-\frac{1}{2}\|u_0\|^2} du_0 \right\|^2 \\ &= \left\| \frac{1}{\mu\kappa} \int_{\mathbb{R}^n} ([f(R_x(\mu u)) - f(x)] - [f(R_x(\mu u)) - f(x) - \langle \text{grad}f(x), \mu u \rangle]) u e^{-\frac{1}{2}\|u_0\|^2} du_0 \right\|^2 \\ &\leq 2\|\mathbb{E}(g_\mu(x))\|^2 + \left\| \frac{2}{\mu^2} \int_{\mathbb{R}^n} (f(R_x(\mu u)) - f(x) - \langle \text{grad}f(x), \mu u \rangle) u e^{-\frac{1}{2}\|u_0\|^2} du_0 \right\|^2 \\ &\leq 2\|\mathbb{E}(g_\mu(x))\|^2 + \frac{2}{\mu^2} \int_{\mathbb{R}^n} (f(R_x(\mu u)) - f(x) - \langle \text{grad}f(x), \mu u \rangle)^2 \|u\|^2 e^{-\frac{1}{2}\|u_0\|^2} du_0 \\ &\leq 2\|\mathbb{E}(g_\mu(x))\|^2 + \frac{\mu^2}{2} L_g (d+6)^3, \end{aligned}$$

where the last inequality is from the same trick as in part (a). This completes the proof of part (b).

Finally, we prove part (c). Since  $\mathbb{E}(\|g_\mu(x)\|^2) = \frac{1}{\mu^2} \mathbb{E}_{u_0} [(f(R_x(\mu u)) - f(x))^2 \|u\|^2]$ , and  $(f(R_x(\mu u)) - f(x))^2 = (f(R_x(\mu u)) - f(x) - \mu \langle \text{grad}f(x), u \rangle + \mu \langle \text{grad}f(x), u \rangle)^2 \leq 2(\frac{L_g}{2} \mu^2 \|u\|^2)^2 + 2\mu^2 \langle \text{grad}f(x), u \rangle^2$ , we have

$$\mathbb{E}(\|g_\mu(x)\|^2) \leq \frac{\mu^2}{2} L_g^2 \mathbb{E}(\|u\|^6) + 2\mathbb{E}(\|\langle \text{grad}f(x), u \rangle u\|^2) \leq \frac{\mu^2}{2} L_g^2 (d+6)^3 + 2\mathbb{E}(\|\langle \text{grad}f(x), u \rangle u\|^2). \quad (2.8)$$

Now we bound the term  $\mathbb{E}(\|\langle \text{grad}f(x), u \rangle u\|^2)$  using the same trick as in [NS17]. Without loss of generality, assume  $\mathcal{X}$  is the  $d$ -dimensional subspace generated by the first  $d$  coordinates, i.e.,  $\forall x \in \mathcal{X}$ , the last  $n-d$  elements of  $x$  are zeros. Also for brevity, denote  $g = \text{grad}f(x)$ . We have that

$$\begin{aligned} \mathbb{E}(\|\langle \text{grad}f(x), u \rangle u\|^2) &= \frac{1}{\kappa} \int_{\mathbb{R}^n} \langle \text{grad}f(x), u \rangle^2 \|u\|^2 e^{-\frac{1}{2}\|u_0\|^2} du_0 \\ &= \frac{1}{\kappa(d)} \int_{\mathbb{R}^d} \left( \sum_{i=1}^d g_i x_i \right)^2 \left( \sum_{i=1}^d x_i^2 \right) e^{-\frac{1}{2} \sum_{i=1}^d x_i^2} dx_1 \cdots dx_d, \end{aligned}$$

where  $x_i$  denotes the  $i$ -th coordinate of  $u_0$ , the last  $n-d$  dimensions are integrated to be one, and  $\kappa(d)$  is the normalization constant for  $d$ -dimensional Gaussian distribution. For simplicity, denote

$x = (x_1, \dots, x_d)$ , then

$$\begin{aligned}
\mathbb{E}(\|\langle \text{grad}f(x), u \rangle u\|^2) &= \frac{1}{\kappa(d)} \int_{\mathbb{R}^d} \langle g, x \rangle^2 \|x\|^2 e^{-\frac{1}{2}\|x\|^2} dx \\
&\leq \frac{1}{\kappa(d)} \int_{\mathbb{R}^d} \|x\|^2 e^{-\frac{\tau}{2}\|x\|^2} \langle g, x \rangle^2 e^{-\frac{1-\tau}{2}\|x\|^2} dx \leq \frac{2}{\kappa(d)\tau e} \int_{\mathbb{R}^d} \langle g, x \rangle^2 e^{-\frac{1-\tau}{2}\|x\|^2} dx \\
&= \frac{2}{\kappa(d)\tau(1-\tau)^{1+d/2}e} \int_{\mathbb{R}^d} \langle g, x \rangle^2 e^{-\frac{1}{2}\|x\|^2} dx = \frac{2}{\tau(1-\tau)^{1+d/2}e} \|g\|^2,
\end{aligned} \tag{2.9}$$

where the second inequality is due to the following fact:  $x^p e^{-\frac{\tau}{2}x^2} \leq (\frac{2}{\tau e})^{p/2}$ . Taking  $\tau = \frac{2}{(d+4)}$  gives the desired result.  $\square$

### 2.3 The Zeroth-order Riemannian Hessian Estimator

We now extend the above methodology and propose estimators for the Riemannian Hessian in the stochastic zeroth-order setting. We restrict our discussion to compact submanifolds embedded in Euclidean space, so that the definition of Riemannian Hessian (2.4) is applied. We assume the following assumption of  $F(x, \xi)$ :

**Assumption 2.2.** *Given any point  $x \in \mathcal{M}$  and  $\eta \in T_x\mathcal{M}$ , we have*

$$\|P_\eta^{-1} \circ \text{Hess}F(R_x(\eta), \xi) \circ P_\eta - \text{Hess}F(x, \xi)\|_{\text{op}} \leq L_H \|\eta\|, \tag{2.10}$$

almost everywhere for  $\xi$ , where  $P_\eta : T_x\mathcal{M} \rightarrow T_{R_x(\eta)}\mathcal{M}$  denotes the parallel transport [ABBC20], an isometry from the tangent space of  $x$  to the tangent space of  $R_x(\eta)$ , and  $\circ$  is the function composition. Here  $\|\cdot\|_{\text{op}}$  is the operator norm in the ambient Euclidean space.

Assumption 2.2 is the analogue of the Lipschitz Hessian type assumption from the Euclidean setting, and induces the following equivalent conditions (see, also [ABBC20]):

$$\begin{aligned}
\|P_\eta^{-1} \text{grad}F(R_x(\eta), \xi) - \text{grad}f(x) - \text{Hess}F(x, \xi)[\eta]\| &\leq \frac{L_H}{2} \|\eta\|^2 \\
\left| F(R_x(\eta), \xi) - \left[ F(x, \xi) + \langle \eta, \text{grad}F(x, \xi) \rangle + \frac{1}{2} \langle \eta, \text{Hess}F(x, \xi)[\eta] \rangle \right] \right| &\leq \frac{L_H}{6} \|\eta\|^3.
\end{aligned} \tag{2.11}$$

In the Euclidean setting,  $P_\eta$  reduces to the identity mapping. Throughout this section, we also assume that  $F(\cdot, \xi)$  satisfies Assumption 2.1 and the following assumption, which is used frequently in zeroth-order stochastic optimization [GL13, BG19, ZYYF19].

**Assumption 2.3.** *We have (with  $\mathbb{E} = \mathbb{E}_\xi$ ) that,  $\mathbb{E}[F(x, \xi)] = f(x)$ ,  $\mathbb{E}[\text{grad}F(x, \xi)] = \text{grad}f(x)$  and  $\mathbb{E}[\|\text{grad}F(x, \xi) - \text{grad}f(x)\|^2] \leq \sigma^2$ ,  $\forall x \in \mathcal{M}$ .*

We first introduce the following identity which follows immediately from the second-order Stein's identity for Gaussian distribution [Ste72].

**Lemma 2.3.** *Suppose  $\mathcal{X}$  is a  $d$ -dimensional subspace of  $\mathbb{R}^n$ , with orthogonal projection matrix  $P \in \mathbb{R}^{n \times n}$ ,  $P = P^2 = P^\top$ , and  $u_0 \sim \mathcal{N}(0, I_n)$  is a standard normal distribution and  $u = Pu_0$  is the orthogonal projection of  $u_0$  onto the subspace. Then  $\forall H \in \mathbb{R}^{n \times n}$ ,  $H^\top = H$ , and  $H = PHP$  (which means that the eigenvectors of  $H$  lies all in  $\mathcal{X}$ ), we have*

$$PHP = \frac{1}{2\kappa} \int_{\mathbb{R}^n} \langle u, Hu \rangle (uu^\top - P) e^{-\frac{1}{2}\|u_0\|^2} du_0 = \mathbb{E} \left[ \frac{1}{2} \langle u, Hu \rangle (uu^\top - P) \right], \tag{2.12}$$

where  $\|\cdot\|$  here is the Euclidean norm on  $\mathbb{R}^n$ , and  $\kappa$  is the constant for normal density function given by  $\kappa := \int_{\mathbb{R}^n} e^{-\frac{1}{2}\|u\|^2} du = (2\pi)^{n/2}$ .

The identity in (2.12) simply follows by applying the second-order Stein's identity,  $\mathbb{E}[(xx^\top - I_n)g(x)] = \mathbb{E}[\nabla^2 g(x)]$ , directly to the function  $g(x) = \frac{1}{2}\langle x, Hx \rangle$  and multiplying the resulting identity by  $P$  on both sides.

**Lemma 2.4.** [BG19] *Suppose  $\mathcal{X}$  is a  $d$ -dimensional subspace of  $\mathbb{R}^n$ , with orthogonal projection matrix  $P \in \mathbb{R}^{n \times n}$ ,  $P = P^2 = P^\top$ , and  $u_0 \sim \mathcal{N}(0, I_n)$  is a standard normal distribution and  $u = Pu_0$  is the orthogonal projection of  $u_0$  onto the subspace. Then*

$$\mathbb{E}[\|u_0 u_0^\top - I_n\|_F^8] \leq 2(n+16)^8 \quad \text{and} \quad \mathbb{E}[\|uu^\top - P\|_F^8] \leq 2(d+16)^8. \quad (2.13)$$

*Proof.* Proof of Lemma 2.4 See [BG19] for the proof of the first inequality in Eq. (2.13). We now show how to get the right part from the left. Similar to the proof of Corollary 2.1, we use an eigen-decomposition of  $P = Q^\top \Lambda Q$  and get (again  $\tilde{u} = Qu$ ):

$$\mathbb{E}\|uu^\top - P\|_F^8 = \mathbb{E}\|(\tilde{u}_1, \dots, \tilde{u}_d)^\top (\tilde{u}_1, \dots, \tilde{u}_d) - I_d\|_F^8 \leq 2(d+16)^8,$$

which completes the proof.  $\square$

We now propose our zeroth-order Riemannian Hessian estimator, motivated by the zeroth-order Hessian estimator in the Euclidean setting proposed by [BG19].

**Definition 2.7** (Zeroth-Order Riemannian Hessian). *Generate  $u \in T_x \mathcal{M}$  following a standard normal distribution on the tangent space  $T_x \mathcal{M}$ , by projection  $u = P_x u_0$  as described in Section 2.2. Then, the zeroth-order Riemannian Hessian estimator of a function  $f$  at the point  $x$  is given by*

$$H_\mu(x) = \frac{1}{2\mu^2} (uu^\top - P) [F(R_x(\mu u), \xi) + F(R_x(-\mu u), \xi) - 2F(x, \xi)]. \quad (2.14)$$

*Note that our Riemannian Hessian estimator is actually the Hessian estimator of the pullback function  $\hat{F}_x(\eta, \xi) = F(R_x(\eta), \xi)$ ,  $\forall x \in \mathcal{M}$  and  $\eta \in T_x \mathcal{M}$  projected onto the tangent space  $T_x \mathcal{M}$ .*

We immediately have the following bound on the variance of  $H_\mu(x)$ .

**Lemma 2.5.** *Under Assumption 2.1, the Riemannian Hessian estimator given in Eq. (2.14) satisfies*

$$\mathbb{E}_{u, \Xi} \|H_\mu(x)\|_F^4 \leq \frac{(d+16)^8}{8} L_g^2. \quad (2.15)$$

*Proof.* Proof of Lemma 2.5 From Assumption 2.1 and Corollary 2.1 we have

$$\begin{aligned} & \mathbb{E}|F(R_x(\mu u), \xi) + F(R_x(-\mu u), \xi) - 2F(x, \xi)|^8 \\ &= \mathbb{E}|F(R_x(\mu u), \xi) - F(x, \xi) - \langle \text{grad} F(x, \xi), \mu u \rangle + F(R_x(-\mu u), \xi) - F(x, \xi) - \langle \text{grad} F(x, \xi), -\mu u \rangle|^8 \\ &\leq \mathbb{E}\left[\frac{\mu^2 L_g}{2} \|u\|^2 + \frac{\mu^2 L_g}{2} \|u\|^2\right]^8 = \mathbb{E}[\mu^{16} L_g^8 \|u\|^{16}] \leq \mu^{16} L_g^8 (d+16)^8. \end{aligned} \quad (2.16)$$

Moreover, we have

$$\begin{aligned} \mathbb{E}\|H_\mu(x)\|_F^4 &= \mathbb{E}\left\| \frac{1}{2\mu^2} (uu^\top - P) [F(R_x(\mu u), \xi) + F(R_x(-\mu u), \xi) - 2F(x, \xi)] \right\|_F^4 \\ &\leq \frac{1}{16\mu^8} \left( \mathbb{E}|F(R_x(\mu u), \xi) + F(R_x(-\mu u), \xi) - 2F(x, \xi)|^8 \mathbb{E}\|uu^\top - P\|_F^8 \right)^{1/2} \\ &\leq \frac{(d+16)^4}{8\mu^8} \left( \mathbb{E}|F(R_x(\mu u), \xi) + F(R_x(-\mu u), \xi) - 2F(x, \xi)|^8 \right)^{1/2}, \end{aligned} \quad (2.17)$$

where the first inequality is by Hölder's inequality and the second one is by Lemma 2.4. Combining (2.16) and (2.17) yields the desired result (2.15).  $\square$

We will also use the mini-batch multi-sampling technique. For  $i = 1, \dots, b$ , denote each Hessian estimator as

$$H_{\mu,i}(x) = \frac{1}{2\mu^2}(u_i u_i^\top - P)[F(R_x(\mu u_i), \xi_i) + F(R_x(-\mu u_i), \xi_i) - 2F(x, \xi_i)]. \quad (2.18)$$

The averaged Hessian estimator is given by

$$\bar{H}_{\mu,\xi}(x) = \frac{1}{b} \sum_{i=1}^b H_{\mu,i}(x). \quad (2.19)$$

We now have the following bound of  $\bar{H}_{\mu,\xi}(x)$  and  $\text{Hess}f(x)$ .

**Lemma 2.6.** *Under Assumption 2.1 and Assumption 2.2, let  $\bar{H}_{\mu,\xi}(x)$  be calculated as in Eq. (2.19), then we have that:  $\forall x \in \mathcal{M}$  and  $\forall \eta \in T_x \mathcal{M}$ ,*

$$\mathbb{E}_{\mathcal{U},\Xi} \|\bar{H}_{\mu,\xi}(x) - \text{Hess}f(x)\|_{\text{op}}^2 \leq \frac{(d+16)^4}{\sqrt{2b}} L_g + \frac{\mu^2 L_H^2}{18} (d+6)^5, \quad (2.20)$$

$$\mathbb{E}_{\mathcal{U},\Xi} \|\bar{H}_{\mu,\xi}(x) - \text{Hess}f(x)\|_{\text{op}}^3 \leq \tilde{C} \frac{(d+16)^6}{b^{3/2}} L_g^{1.5} + \frac{1}{27} \mu^3 L_H^3 (d+6)^{7.5}, \quad (2.21)$$

where  $\|\cdot\|_{\text{op}}$  denotes the operator norm and  $\tilde{C}$  is some absolute constant.

*Proof.* Proof of Lemma 2.6 Denote  $\mathbb{E} = \mathbb{E}_{\mathcal{U},\Xi}$  as the expectation with respect to all previous random variables. We first show Eq. (2.20). Denote  $X_i = H_{\mu,i} - \mathbb{E}H_{\mu,i}$ , then  $X_i$ 's are iid zero-mean random matrices. Since  $\|\cdot\|_{\text{op}} \leq \|\cdot\|_F$ , we have

$$\begin{aligned} \mathbb{E} \|\bar{H}_{\mu,\xi}(x) - \mathbb{E}\bar{H}_{\mu,\xi}(x)\|_{\text{op}}^2 &= \mathbb{E} \left\| \frac{1}{b} \sum_{i=1}^b X_i \right\|_{\text{op}}^2 \leq \mathbb{E} \left\| \frac{1}{b} \sum_{i=1}^b X_i \right\|_F^2 \\ &= \mathbb{E} \left[ \frac{1}{b^2} \sum_{i=1}^b \|X_i\|_F^2 + \frac{1}{b^2} \sum_{i \neq j} \langle X_i, X_j \rangle \right] = \mathbb{E} \left[ \frac{1}{b^2} \sum_{i=1}^b \|X_i\|_F^2 \right] \\ &= \mathbb{E} \frac{1}{b^2} b \|X_1\|_F^2 = \mathbb{E} \frac{1}{b} \|H_{\mu,1} - \mathbb{E}H_{\mu,1}\|_F^2 = \frac{1}{b} \mathbb{E} [\|H_{\mu,1}\|_F^2 - \|\mathbb{E}H_{\mu,1}\|_F^2] \\ &\leq \frac{1}{b} \mathbb{E} \|H_{\mu,1}\|_F^2 \leq \frac{1}{b} \sqrt{\mathbb{E} \|H_{\mu,1}(x)\|_F^4} \leq \frac{(d+16)^4}{2\sqrt{2b}} L_g, \end{aligned} \quad (2.22)$$

where the third inequality is from the Jensen's inequality, and the last inequality is due to Eq. (2.15). Note that (2.22) immediately implies

$$\begin{aligned} \mathbb{E} \|\bar{H}_{\mu,\xi}(x) - \text{Hess}f(x)\|_{\text{op}}^2 &\leq 2\mathbb{E} \|\bar{H}_{\mu,\xi}(x) - \mathbb{E}\bar{H}_{\mu,\xi}(x)\|_{\text{op}}^2 + 2\|\mathbb{E}\bar{H}_{\mu,\xi}(x) - \text{Hess}f(x)\|_{\text{op}}^2 \\ &\leq \frac{(d+16)^4}{\sqrt{2b}} L_g + 2\|\mathbb{E}\bar{H}_{\mu,\xi}(x) - \text{Hess}f(x)\|_{\text{op}}^2. \end{aligned} \quad (2.23)$$

Now we bound the term  $\|\mathbb{E}\bar{H}_{\mu,\xi}(x) - \text{Hess}f(x)\|_{\text{op}}^2$ . Note that

$$\begin{aligned}
& |\langle \eta, (\mathbb{E}H_{\mu,i}(x) - \text{Hess}f(x))[\eta] \rangle| \\
&= \left| \langle \eta, \left( \mathbb{E} \left[ \frac{1}{2\mu^2} (uu^\top - P)[f(R_x(\mu u)) + f(R_x(-\mu u)) - 2f(x)] \right] - \text{Hess}f(x) \right) [\eta] \right\rangle \right| \\
&= \left| \langle \eta, \left( \mathbb{E} \left[ \frac{1}{2\mu^2} (uu^\top - P)[f(R_x(\mu u)) + f(R_x(-\mu u)) - 2f(x) - \mu^2 \langle u, \text{Hess}f(x)[u] \rangle] \right) [\eta] \right\rangle \right| \\
&= \frac{1}{2\mu^2} \left| \langle \eta, \left( \mathbb{E} \left[ [f(R_x(\mu u)) - f(x) - \frac{\mu^2}{2} \langle u, \text{Hess}f(x)[u] \rangle \right. \right. \right. \\
&\quad \left. \left. \left. + f(R_x(-\mu u)) - f(x) - \frac{\mu^2}{2} \langle u, \text{Hess}f(x)[u] \rangle](uu^\top - P) \right] \right) [\eta] \right\rangle \right|,
\end{aligned}$$

which together with Assumption 2.2 yields

$$\begin{aligned}
& |\langle \eta, (\mathbb{E}H_{\mu,i}(x) - \text{Hess}f(x))[\eta] \rangle| \leq \frac{\mu L_H}{6} \mathbb{E} \left[ \|u\|^3 \|uu^\top - P\|_{\text{op}} \right] \|\eta\|^2 \\
(\text{H\"older}) & \leq \frac{\mu L_H}{6} \sqrt{\mathbb{E}\|u\|^6 \mathbb{E}\|uu^\top - P\|_F^2} \|\eta\|^2 \leq \frac{\mu L_H}{6} (d+6)^{5/2} \|\eta\|^2,
\end{aligned} \tag{2.24}$$

where the last inequality is by Corollary 2.1 and Lemma 2.4. (2.24) implies

$$\|\mathbb{E}\bar{H}_{\mu,\xi}(x) - \text{Hess}f(x)\|_{\text{op}} \leq \frac{\mu L_H}{6} (d+6)^{5/2}. \tag{2.25}$$

Combining (2.23) and (2.25) gives Eq. (2.20).

Now we show Eq. (2.21). By a similar analysis we have

$$\begin{aligned}
& \mathbb{E}\|\bar{H}_{\mu,\xi}(x) - \text{Hess}f(x)\|_{\text{op}}^3 \\
& \leq \mathbb{E}(\|\bar{H}_{\mu,\xi}(x) - \mathbb{E}\bar{H}_{\mu,\xi}(x)\|_{\text{op}} + \|\mathbb{E}\bar{H}_{\mu,\xi}(x) - \text{Hess}f(x)\|_{\text{op}})^3 \\
& \leq 8\mathbb{E}\|\bar{H}_{\mu,\xi}(x) - \mathbb{E}\bar{H}_{\mu,\xi}(x)\|_{\text{op}}^3 + 8\|\mathbb{E}\bar{H}_{\mu,\xi}(x) - \text{Hess}f(x)\|_{\text{op}}^3 \\
(\text{H\"older}) & \leq 8\sqrt{\mathbb{E}\|\bar{H}_{\mu,\xi}(x) - \mathbb{E}\bar{H}_{\mu,\xi}(x)\|_{\text{op}}^2 \mathbb{E}\|\bar{H}_{\mu,\xi}(x) - \mathbb{E}\bar{H}_{\mu,\xi}(x)\|_{\text{op}}^4} \\
& \quad + 8\|\mathbb{E}\bar{H}_{\mu,\xi}(x) - \text{Hess}f(x)\|_{\text{op}}^3,
\end{aligned} \tag{2.26}$$

where the second inequality is by the following fact: when  $a, b \geq 0$ ,  $(a+b)^3 \leq \max\{(2a)^3, (2b)^3\} \leq 8a^3 + 8b^3$ . Moreover, since  $\|\cdot\|_{\text{op}} \leq \|\cdot\|_F$ , and  $X_i = H_{\mu,i} - \mathbb{E}H_{\mu,i}$  are iid zero-mean random matrices,

we have

$$\begin{aligned}
& \mathbb{E}\|\bar{H}_{\mu,\xi}(x) - \mathbb{E}\bar{H}_{\mu,\xi}(x)\|_{\text{op}}^4 = \mathbb{E}\left\|\frac{1}{b}\sum_{i=1}^b X_i\right\|_{\text{op}}^4 \leq \frac{C}{b^4} \left( \mathbb{E}\left\|\sum_{i=1}^b X_i\right\|_{\text{op}} + (b\mathbb{E}\|X_i\|_{\text{op}}^4)^{1/4} \right)^4 \\
& \leq \frac{C}{b^4} \left( \sqrt{\mathbb{E}\left\|\sum_{i=1}^b X_i\right\|_F^2} + (b\mathbb{E}\|X_i\|_F^4)^{1/4} \right)^4 = \frac{C}{b^4} \left( \sqrt{\sum_{i=1}^b \mathbb{E}\|X_i\|_F^2} + (b\mathbb{E}\|X_i\|_F^4)^{1/4} \right)^4 \\
& = \frac{C}{b^4} \left( \sqrt{b}\sqrt{\mathbb{E}\|X_1\|_F^2} + (b\mathbb{E}\|X_1\|_F^4)^{1/4} \right)^4 \leq \frac{C}{b^4} \left( \sqrt{b}\sqrt[4]{\mathbb{E}\|X_1\|_F^4} + (b\mathbb{E}\|X_1\|_F^4)^{1/4} \right)^4 \\
& = \frac{C}{b^4} (\sqrt{b} + \sqrt[4]{b})^4 \mathbb{E}\|H_{\mu,1} - \mathbb{E}H_{\mu,1}\|_F^4 \leq \frac{16C}{b^2} \mathbb{E}\|H_{\mu,1} - \mathbb{E}H_{\mu,1}\|_F^4 \tag{2.27} \\
& = \frac{16C}{b^2} \mathbb{E}(\|H_{\mu,1}\|_F^2 - 2\langle H_{\mu,1}, \mathbb{E}H_{\mu,1} \rangle + \|\mathbb{E}H_{\mu,1}\|_F^2)^2 \\
& \leq \frac{16C}{b^2} \mathbb{E}(\|H_{\mu,1}\|_F^2 + 2\|H_{\mu,1}\|_F \|\mathbb{E}H_{\mu,1}\|_F + \|\mathbb{E}H_{\mu,1}\|_F^2)^2 \\
& \leq \frac{16C}{b^2} \mathbb{E}(2\|H_{\mu,1}\|_F^2 + 2\|\mathbb{E}H_{\mu,1}\|_F^2)^2 \leq \frac{16C}{b^2} \mathbb{E}(2\|H_{\mu,1}\|_F^2 + 2\mathbb{E}\|H_{\mu,1}\|_F^2)^2 \\
& \leq \frac{64C}{b^2} (\mathbb{E}\|H_{\mu,1}\|_F^4 + \mathbb{E}\|H_{\mu,1}\|_F^4) \leq \frac{128C}{b^2} (d+16)^8 L_g^2,
\end{aligned}$$

where the first inequality is due to the Rosenthal inequality [Rio09],  $C$  is an absolute constant, the fourth inequality is due to the fact  $1 \leq \sqrt[4]{b} \leq \sqrt{b}$ . Plugging Eq. (2.22), Eq. (2.25) and Eq. (2.27) back to Eq. (2.26) gives the desired result (2.21).  $\square$

### 3 Stochastic Zeroth-order Riemannian Optimization Algorithms

We now demonstrate the applicability of the developed Riemannian derivative estimation methodology in Section 2, for various classes of stochastic zeroth-order Riemannian optimization algorithms.

#### 3.1 Zeroth-order Smooth Riemannian Optimization

In this section, we focus on the smooth optimization problem with  $h \equiv 0$  and  $f$  satisfying Assumption 2.1. We propose Z0-RGD, the zeroth-order Riemannian gradient descent method and provide its complexity analysis. The algorithm is formally presented in Algorithm 1.

---

#### Algorithm 1 Zeroth-Order Riemannian Gradient Descent (Z0-RGD)

---

- 1: **Input:** Initial point  $x_0 \in \mathcal{M}$ , smoothing parameter  $\mu$ , step size  $\eta_k$ , fixed number of iteration  $N$ .
  - 2: **for**  $k = 0$  **to**  $N - 1$  **do**
  - 3:   Sample a standard Gaussian random vector  $u_k \in T_{x_k}\mathcal{M}$  by orthogonal projection in Definition 2.6.
  - 4:   Compute the zeroth-order gradient  $g_\mu(x_k)$  by Eq. (2.5).
  - 5:   Update  $x_{k+1} = R_{x_k}(-\eta_k g_\mu(x_k))$ .
  - 6: **end for**
- 

The following theorem gives the iteration and oracle complexities of Algorithm 1 for obtaining an  $\epsilon$ -stationary point of (1.1) when  $h \equiv 0$ .

**Theorem 3.1.** *Let  $f$  satisfy Assumption 2.1 and suppose  $\{x_k\}$  is the sequence generated by Algorithm 1 with the stepsize  $\eta_k = \hat{\eta} = \frac{1}{2(d+4)L_g}$ . Then, we have*

$$\frac{1}{N+1} \sum_{k=0}^N \mathbb{E}_{\mathcal{U}_k} \|\text{grad}f(x_k)\|^2 \leq \frac{4}{\hat{\eta}} \left( \frac{f(x_0) - f(x^*)}{N+1} + C(\mu) \right), \quad (3.1)$$

where  $\mathcal{U}_k$  denotes the set of all Gaussian random vectors we drew for the first  $k$  iterations<sup>3</sup>, and  $C(\mu) = \frac{\mu^2 L_g (d+3)^3}{16 (d+4)} + \frac{\mu^2 (d+6)^3}{16 (d+4)} + \frac{\mu^2 L_g (d+6)^3}{16 (d+4)^2}$ . In order to have

$$\frac{1}{N+1} \sum_{k=0}^N \mathbb{E}_{\mathcal{U}_k} \|\text{grad}f(x_k)\|^2 \leq \epsilon^2, \quad (3.2)$$

we need the smoothing parameter  $\mu$  and number of iteration  $N$  (which is also the number of calls to the zeroth-order oracle) to be set as  $\mu = \mathcal{O}(\epsilon/d^{3/2})$ ,  $N = \mathcal{O}(d/\epsilon^2)$ .

*Proof.* Proof of Theorem 3.1 From Assumption 2.1 we have

$$f(x_{k+1}) \leq f(x_k) - \eta_k \langle g_\mu(x_k), \text{grad}f(x_k) \rangle + \frac{\eta_k^2 L_g}{2} \|g_\mu(x_k)\|^2.$$

Taking the expectation w.r.t.  $u_k$  on both sides, we have

$$\begin{aligned} \mathbb{E}_{u_k} [f(x_{k+1})] &\leq f(x_k) - \eta_k \langle \mathbb{E}_{u_k}(g_\mu(x_k)), \text{grad}f(x_k) \rangle + \frac{\eta_k^2 L_g}{2} \mathbb{E}_{u_k} (\|g_\mu(x_k)\|^2) \\ &\leq f(x_k) - \eta_k \langle \mathbb{E}_{u_k}(g_\mu(x_k)), \text{grad}f(x_k) \rangle + \frac{\eta_k^2 L_g}{2} \left( \frac{\mu^2}{2} L_g^2 (d+6)^3 + 2(d+4) \|\text{grad}f(x_k)\|^2 \right), \end{aligned}$$

where the last inequality is by Proposition 2.1. Now Take  $\eta_k = \hat{\eta} = \frac{1}{2(d+4)L_g}$ , we have

$$\begin{aligned} &\mathbb{E}_{u_k} [f(x_{k+1})] \\ &\leq f(x_k) + \frac{\hat{\eta}}{2} (\|\text{grad}f(x_k)\|^2 - 2 \langle \mathbb{E}_{u_k}(g_\mu(x_k)), \text{grad}f(x_k) \rangle) + \frac{\mu^2 L_g (d+6)^3}{16 (d+4)^2} \\ &= f(x_k) + \frac{\hat{\eta}}{2} (\|\text{grad}f(x_k) - \mathbb{E}_{u_k}(g_\mu(x_k))\|^2 - \|\mathbb{E}_{u_k}(g_\mu(x_k))\|^2) + \frac{\mu^2 L_g (d+6)^3}{16 (d+4)^2} \\ &\leq f(x_k) + \frac{\hat{\eta}}{2} \left( \frac{\mu^2 L_g^2}{4} (d+3)^3 - \frac{1}{2} \|\text{grad}f(x_k)\|^2 + \frac{\mu^2}{4} L_g (d+6)^3 \right) + \frac{\mu^2 L_g (d+6)^3}{16 (d+4)^2} \\ &= f(x_k) - \frac{\hat{\eta}}{4} \|\text{grad}f(x_k)\|^2 + C(\mu), \end{aligned}$$

where the second inequality is from Proposition 2.1. Define  $\phi_k := f(x_k) - f(x^*)$ . Now take the expectation w.r.t.  $\mathcal{U}_k = \{u_0, u_1, \dots, u_{k-1}\}$ , we have

$$\phi_{k+1} \leq \phi_k - \frac{\hat{\eta}}{4} \mathbb{E}_{\mathcal{U}_k} \|\text{grad}f(x_k)\|^2 + C(\mu).$$

Summing the above inequality over  $k = 0, \dots, N$  yields (3.1).

Therefore with  $\mu = \mathcal{O}(\epsilon/d^{3/2})$  we have  $C(\mu) \leq \hat{\eta} \epsilon^2 / 4$ . Taking  $N \geq 8(d+4)L_g(f(x_0) - f(x^*)) / \epsilon^2$  yields (3.2). In summary, the number of iterations for obtaining an  $\epsilon$ -stationary solution is  $\mathcal{O}(d/\epsilon^2)$ , and hence the total zeroth-order oracle complexity is also  $\mathcal{O}(d/\epsilon^2)$ .  $\square$

<sup>3</sup>The notation of taking the expectation w.r.t. a set, is to take the expectation for each of the elements in the set.

**Remark 3.1.** Note that in Algorithm 1, we only sample one Gaussian vector in each iteration of the algorithm. In practice, one can also sample multiple Gaussian random vectors in each iteration and obtain an averaged gradient estimator. Suppose we sample  $m$  i.i.d. Gaussian random vectors in each iteration and use the average  $\bar{g}_\mu(x) = \frac{1}{m} \sum_{i=1}^m g_{\mu,i}(x)$ , then the bound for our zeroth-order estimator becomes

$$\mathbb{E}(\|\bar{g}_\mu(x) - \text{grad}f(x)\|^2) \leq \mu^2 L_g^2 (d+6)^3 + \frac{2(d+4)}{m} \|\text{grad}f(x)\|^2. \quad (3.3)$$

Hence, the final result in Theorem 3.1 can be improved to

$$\frac{1}{N+1} \sum_{k=0}^N \mathbb{E}_{\mathcal{U}_k} \|\text{grad}f(x_k)\|^2 \leq 4L_g \frac{f(x_0) - f(x^*)}{N+1} + \mu^2 L_g^2 (d+6)^3, \quad (3.4)$$

with  $\hat{\eta} = 1/L_g$  and  $C(\mu) = \mu^2 L_g^2 (d+6)^3 / 2$ . Therefore the number of iterations required is improved to  $N = \mathcal{O}(1/\epsilon^2)$  when we set  $\mu = \mathcal{O}(\epsilon/d^{3/2})$  and  $m = \mathcal{O}(d)$ . However, the zeroth-order oracle complexity is still  $\mathcal{O}(d/\epsilon^2)$ . The proof of (3.3) and (3.4) is given in the appendix. This multi-sampling technique will play a key role in our stochastic and non-smooth case analyses.

### 3.2 Zeroth-Order Stochastic Riemannian Optimization for Nonconvex Problem

In this section, we focus on the following nonconvex smooth problem:

$$\min_{x \in \mathcal{M}} f(x) := \int_{\xi} F(x, \xi) dP(\xi), \quad (3.5)$$

where  $P$  is a random distribution,  $F$  is a function satisfying Assumption 2.1, in variable  $x$ , almost surely. Note that  $f$  automatically satisfies Assumption 2.1 by the Jensen's inequality.

In the stochastic case, sampling multiple times in every iteration can improve the convergence rate. Our zeroth-order Riemannian gradient estimator is given by

$$\bar{g}_{\mu,\xi}(x) = \frac{1}{m} \sum_{i=1}^m g_{\mu,\xi_i}(x), \text{ where } g_{\mu,\xi_i}(x) = \frac{F(R_x(\mu u_i), \xi_i) - F(x, \xi_i)}{\mu} u_i, \quad (3.6)$$

and  $u_i$  is a standard normal random vector on  $T_x \mathcal{M}$ . We also immediately have that

$$\mathbb{E}_{\xi_i} g_{\mu,\xi_i}(x) = \frac{f(R_x(\mu u)) - f(x)}{\mu} u = g_\mu(x). \quad (3.7)$$

The multi-sampling technique enables us to obtain the following bound on  $\mathbb{E}\|\bar{g}_{\mu,\xi}(x) - \text{grad}f(x)\|^2$ , the proof of which is given in the Appendix C.

**Lemma 3.1.** For the Riemannian gradient estimator in (3.6), under Assumptions 2.1 and 2.3, we have

$$\mathbb{E}\|\bar{g}_{\mu,\xi}(x) - \text{grad}f(x)\|^2 \leq \mu^2 L_g^2 (d+6)^3 + \frac{8(d+4)}{m} \sigma^2 + \frac{8(d+4)}{m} \|\text{grad}f(x)\|^2, \quad (3.8)$$

where the expectation  $\mathbb{E}$  is taken for both Gaussian vectors  $\mathcal{U} = \{u_1, \dots, u_m\}$  and  $\xi$ .

Our zeroth-order Riemannian stochastic gradient descent algorithm (ZO-RSGD) for solving (3.5), is presented in Algorithm 2.

Now we present convergence analysis for obtaining an  $\epsilon$ -stationary point of (3.5).



---

**Algorithm 2** Zeroth-order Riemannian Stochastic Gradient Descent (ZO-RSGD)

---

- 1: **Input:** Initial point  $x_0 \in \mathcal{M}$ , smoothing parameter  $\mu$ , multi-sample constant  $m$ , step size  $\eta_k$ , fixed number of iteration  $N$ .
  - 2: **for**  $k = 0$  **to**  $N - 1$  **do**
  - 3:   Sample the standard Gaussian random vectors  $u_i^k$  on  $T_{x_k}\mathcal{M}$  by orthogonal projection in Definition 2.6, and sample  $\xi_i^k$ ,  $i = 1, \dots, m$ .
  - 4:   Compute the zeroth-order gradient  $\bar{g}_{\mu, \xi}(x_k)$  by Eq. (3.6).
  - 5:   Update  $x_{k+1} = R_{x_k}(-\eta_k \bar{g}_{\mu, \xi}(x_k))$ .
  - 6: **end for**
- 

**Theorem 3.2.** *Let  $F$  satisfy Assumption 2.1, w.r.t. variable  $x$  almost surely. Suppose  $\{x_k\}$  is the sequence generated by Algorithm 2 with the stepsize  $\eta_k = \hat{\eta} = \frac{1}{L_g}$ . Under Assumption 2.3, we have*

$$\frac{1}{N+1} \sum_{k=0}^N \mathbb{E}_{\mathcal{U}_k, \Xi_k} \|\text{grad}f(x_k)\|^2 \leq 4L_g \frac{f(x_0) - f(x^*)}{N+1} + C(\mu), \quad (3.9)$$

where  $C(\mu) = 2\mu^2 L_g^2 (d+6)^3 + \frac{16(d+4)}{m} \sigma^2$ ,  $\mathcal{U}_k$  denotes the set of all Gaussian random vectors and  $\Xi_k$  denotes the set of all random variable  $\xi_k$  in the first  $k$  iterations. In order to have  $\frac{1}{N+1} \sum_{k=0}^N \mathbb{E}_{\mathcal{U}_k, \Xi_k} \|\text{grad}f(x_k)\|^2 \leq \epsilon^2$ , we need the smoothing parameter  $\mu$ , number of sampling  $m$  in each iteration and number of iterations  $N$  to be

$$\mu = \mathcal{O}\left(\epsilon/d^{3/2}\right), \quad m = \mathcal{O}\left(d\sigma^2/\epsilon^2\right), \quad N = \mathcal{O}\left(1/\epsilon^2\right). \quad (3.10)$$

Hence, the number of calls to the zeroth-order oracle is  $mN = \mathcal{O}(d/\epsilon^4)$ .

*Proof.* Proof of Theorem 3.2 From Assumption 2.1, we have:

$$f(x_{k+1}) \leq f(x_k) - \eta_k \langle \bar{g}_{\mu, \xi}(x), \text{grad}f(x_k) \rangle + \frac{\eta_k^2 L_g}{2} \|\bar{g}_{\mu, \xi}(x)\|^2$$

Take  $\eta_k = \hat{\eta} = \frac{1}{L_g}$ , we have

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) - \eta_k \langle \bar{g}_{\mu, \xi}(x), \text{grad}f(x_k) \rangle + \frac{\eta_k^2 L_g}{2} \|\bar{g}_{\mu, \xi}(x)\|^2 \\ &= f(x_k) + \frac{1}{2L_g} \left( \|\bar{g}_{\mu, \xi}(x) - \text{grad}f(x)\|^2 - \|\text{grad}f(x)\|^2 \right). \end{aligned}$$

Take the expectation for the random variables at iteration  $k$  on both sides, we have

$$\begin{aligned} \mathbb{E}_k f(x_{k+1}) &\leq f(x_k) + \frac{1}{2L_g} \left( \mathbb{E}_k \|\bar{g}_{\mu, \xi}(x) - \text{grad}f(x)\|^2 - \|\text{grad}f(x)\|^2 \right) \\ \text{Eq. (3.3)} &\leq f(x_k) + \frac{1}{2L_g} \left( \mu^2 L_g^2 (d+6)^3 + \frac{8(d+4)}{m} \sigma^2 + \left( \frac{8(d+4)}{m} - 1 \right) \|\text{grad}f(x)\|^2 \right). \end{aligned}$$

Summing up over  $k = 0, \dots, N$  (assuming that  $m \geq 16(d+4)$ ) yields (3.9). In summary, the total number of iterations for obtaining an  $\epsilon$ -stationary solution of (3.5) is  $\mathcal{O}(1/\epsilon^2)$ , and the stochastic zeroth-order oracle complexity is  $\mathcal{O}(d/\epsilon^4)$ .  $\square$

In Appendix A, we present the oracle complexity of Algorithm 2 when  $f$  is geodesically convex and  $\mathcal{M}$  is the Hadamard manifold.

### 3.3 Zeroth-order Stochastic Riemannian Proximal Gradient Method

We now consider the general optimization problem of the form in Eq. (1.1). For the sake of notation, we denote  $p(x) := f(x) + h(x)$ . We assume that  $\mathcal{M}$  is a compact submanifold,  $h$  is convex in the embedded space  $\mathbb{R}^n$  and is also Lipschitz continuous with parameter  $L_h$ , and  $f(x) := \int_{\xi} F(x, \xi) dP(\xi)$  satisfying Assumption 2.3.

The non-differentiability of  $h$  prohibits Riemannian gradient methods to be applied directly. In [CMMCSZ20], by assuming that the exact gradient of  $f$  is available, a manifold proximal gradient method (ManPG) is proposed for solving (1.1). One typical iteration of ManPG is as follows:

$$\begin{aligned} v_k &:= \operatorname{argmin} \langle \operatorname{grad} f(x_k), v \rangle + \frac{1}{2t} \|v\|^2 + h(x_k + v), \text{ s.t., } v \in T_{x_k} \mathcal{M} \\ x_{k+1} &:= R_{x_k}(\eta_k v_k), \end{aligned} \quad (3.11)$$

where  $t > 0$  and  $\eta_k > 0$  are step sizes. In this section, we develop a zeroth-order counterpart of ManPG (ZO-ManPG), where we assume that only noisy function evaluations of  $f$  are available. The following lemma from [CMMCSZ20] provides a notion of stationary point that is useful for our analysis.

**Lemma 3.2.** *Let  $\bar{v}_k$  be the minimizer of the  $v$ -subproblem in (3.11). If  $\bar{v}_k = 0$ , then  $x_k$  is a stationary point of problem (1.1). We say  $x_k$  is an  $\epsilon$ -stationary point of (1.1) with  $t = \frac{1}{L_g}$ , if  $\|\bar{v}_k\| \leq \epsilon/L_g$ .*

Our ZO-ManPG iterates as:

$$\begin{aligned} v_k &:= \operatorname{argmin} \langle \bar{g}_{\mu, \xi}(x_k), v \rangle + \frac{1}{2t} \|v\|^2 + h(x_k + v), \text{ s.t., } v \in T_{x_k} \mathcal{M}, \\ x_{k+1} &:= R_{x_k}(\eta_k v_k), \end{aligned} \quad (3.12)$$

where  $\bar{g}_{\mu, \xi}(x_k)$  is defined in Eq. (3.6). Note that the only difference between ZO-ManPG (3.12) and ManPG (3.11) is that in (3.12) we use  $\bar{g}_{\mu, \xi}(x)$  to replace the Riemannian gradient  $\operatorname{grad} f$  in (3.11). A more complete description of the algorithm is given in Algorithm 3. Now we provide some useful

---

#### Algorithm 3 Zeroth-Order Riemannian Proximal Gradient Descent (ZO-ManPG)

---

- 1: **Input:** Initial point  $x_0$  on  $\mathcal{M}$ , smoothing parameter  $\mu$ , number of multi-sample  $m$ , step size  $\eta_k$ , fixed number of iteration  $N$ .
  - 2: **for**  $k = 0$  **to**  $N - 1$  **do**
  - 3:   Sample  $m$  standard Gaussian random vector  $u_i$  on  $T_{x_k} \mathcal{M}$  by orthogonal projection in Definition 2.6,  $i = 1, \dots, m$ .
  - 4:   Compute the zeroth-order gradient the random oracle  $\bar{g}_{\mu}(x_k)$  by Eq. (3.6).
  - 5:   Solve  $v_k$  from Eq. (3.12).
  - 6:   Update  $x_{k+1} = R_{x_k}(\eta_k v_k)$ .
  - 7: **end for**
- 

lemmas for analyzing the iteration complexity of Algorithm 3.

**Lemma 3.3.** *(Non-expansiveness) Suppose  $v := \arg \min_{v \in T_x \mathcal{M}} \langle g_1, v \rangle + \frac{1}{2t} \|v\|^2 + h(x + v)$  and  $w := \arg \min_{w \in T_x \mathcal{M}} \langle g_2, w \rangle + \frac{1}{2t} \|w\|^2 + h(x + w)$ . Then we have*

$$\|v - w\| \leq t \|g_1 - g_2\|. \quad (3.13)$$

*Proof.* Proof of Lemma 3.3 By the first order optimality condition [YZS14], we have  $0 \in \frac{1}{t}v + g_1 + \text{Proj}_{T_x\mathcal{M}} \partial h(x+v)$  and  $0 \in \frac{1}{t}w + g_2 + \text{Proj}_{T_x\mathcal{M}} \partial h(x+w)$ , i.e.  $\exists p_1 \in \partial h(x+v)$  and  $p_2 \in \partial h(x+w)$  such that  $v = -t(g_1 + \text{Proj}_{T_x\mathcal{M}}(p_1))$  and  $w = -t(g_2 + \text{Proj}_{T_x\mathcal{M}}(p_2))$ . Therefore we have

$$\begin{aligned}\langle v, w-v \rangle &= t \langle g_1 + \text{Proj}_{T_x\mathcal{M}}(p_1), v-w \rangle \\ \langle w, v-w \rangle &= t \langle g_2 + \text{Proj}_{T_x\mathcal{M}}(p_2), w-v \rangle.\end{aligned}\tag{3.14}$$

Now since  $v, w \in T_x\mathcal{M}$ , and using the convexity of  $h$ , we have

$$\langle \text{Proj}_{T_x\mathcal{M}}(p_1), v-w \rangle = \langle p_1, v-w \rangle = \langle p_1, (v+x) - (w+x) \rangle \geq h(v+x) - h(w+x).\tag{3.15}$$

Substituting Eq. (3.14) and into (3.15) yields,

$$\begin{aligned}\langle v, w-v \rangle &\geq t \langle g_1, v-w \rangle + h(v+x) - h(w+x) \\ \langle w, v-w \rangle &\geq t \langle g_2, w-v \rangle + h(w+x) - h(v+x).\end{aligned}$$

Summing these two inequalities gives  $\langle v-w, v-w \rangle \leq t \langle g_2 - g_1, v-w \rangle$ , and Eq. (3.13) follows by applying the Cauchy-Schwarz inequality.  $\square$

**Corollary 3.1.** *Suppose  $v_k$  is given by (3.12), and  $\bar{v}_k$  is solution of the  $v$ -subproblem in Eq. (3.11), then we have*

$$\mathbb{E}_{\mathcal{U}_k, \Xi_k} \|v_k - \bar{v}_k\|_F^2 \leq t^2 \left( \mu^2 L_g^2 (d+6)^3 + \frac{8(d+4)}{m} \sigma^2 + \frac{8(d+4)}{m} \|\text{grad}f(x_k)\|^2 \right).$$

*Proof.* Proof of Corollary 3.1 By Lemma 3.3, we have

$$\mathbb{E}_{\mathcal{U}_k, \Xi_k} \|v_k - \bar{v}_k\|_F^2 \leq t^2 \mathbb{E}_{\mathcal{U}_k, \Xi_k} \|\bar{g}_{\mu, \xi}(x_k) - \text{grad}f(x_k)\|_F^2.$$

From Lemma 3.1,

$$\begin{aligned}\mathbb{E}_{\mathcal{U}_k, \Xi_k} \|\bar{g}_{\mu, \xi}(x_k) - \text{grad}f(x_k)\|_F^2 \\ \leq \mu^2 L_g^2 (d+6)^3 + \frac{8(d+4)}{m} \sigma^2 + \frac{8(d+4)}{m} \|\text{grad}f(x_k)\|^2.\end{aligned}$$

The desired result hence follows by combining these two inequalities.  $\square$

The following lemma shows the sufficient decrease property for one iteration of **ZO-ManPG**.

**Lemma 3.4.** *For any  $t > 0$ , there exists a constant  $\bar{\eta} > 0$  such that for any  $0 \leq \eta_k \leq \min\{1, \bar{\eta}\}$ , the  $(x_k, v_k)$  generated by Algorithm 3 satisfies*

$$p(x_{k+1}) - p(x_k) \leq - \left( \frac{\eta_k}{2t} - \tilde{C} \right) \|v_k\|^2,\tag{3.16}$$

where  $\tilde{C} = \mu^2 L_g^2 (d+6)^3 + \frac{8(d+4)}{m} \sigma^2 + \frac{8(d+4)}{m} G^2$  and  $G$  is the upper bound of the Riemannian gradient  $\text{grad}f(x)$  (existence by the compactness of  $\mathcal{M}$ ).

*Proof.* Proof of Lemma 3.4 Notice that

$$\begin{aligned}f(x_{k+1}) - f(x_k) &\leq \langle \text{grad}f(x_k), R_{x_k}(\eta_k v_k) - x_k \rangle + \frac{Lg}{2} \|R_{x_k}(\eta_k v_k) - x_k\|^2 \\ &= \langle \text{grad}f(x_k) - \bar{g}_{\mu, \xi}(x), R_{x_k}(\eta_k v_k) - x_k \rangle + \langle \bar{g}_{\mu, \xi}(x), R_{x_k}(\eta_k v_k) - x_k \rangle + \frac{Lg}{2} \|R_{x_k}(\eta_k v_k) - x_k\|^2,\end{aligned}$$

where the inequality follows from Assumption 2.1. Moreover, by Lemma 3.1 and the Fact 3.6 of [CMMCSZ20], we have

$$\begin{aligned} \langle \text{grad}f(x_k) - \bar{g}_{\mu,\xi}(x), R_{x_k}(\eta_k v_k) - x_k \rangle &\leq \|\text{grad}f(x_k) - \bar{g}_{\mu,\xi}(x)\| \|R_{x_k}(\eta_k v_k) - x_k\| \\ &\leq M_1^2 \eta_k^2 \left[ \mu^2 L_g^2 (d+6)^3 + \frac{8(d+4)}{m} \sigma^2 + \frac{8(d+4)}{m} \|\text{grad}f(x)\|^2 \right] \|v_k\|^2. \end{aligned}$$

The rest of the proof of bounding  $\langle \bar{g}_{\mu,\xi}(x), R_{x_k}(\eta_k v_k) - x_k \rangle + \frac{L_g}{2} \|R_{x_k}(\eta_k v_k) - x_k\|^2$  follows from exactly the same process as in ([CMMCSZ20], Lemma 5.2). We omit the details for brevity.  $\square$

**Theorem 3.3.** *Under Assumption 2.3 and Assumption 2.1, the sequence generated by Algorithm 3, with  $\eta_k = \hat{\eta} < \min\{1, \bar{\eta}\}$  and  $t = 1/L_g$ , satisfies:*

$$\frac{1}{N} \sum_{k=0}^{N-1} \mathbb{E}_{\mathcal{U}_k, \Xi_k} \|\bar{v}_k\|^2 \leq \frac{4t(p(x_0) - p(x^*))}{(\hat{\eta} - 8\tilde{C})tN} + \frac{\hat{\eta}Nt^2}{\hat{\eta} - 8\tilde{C}t} \tilde{C} + \frac{8t^3}{\hat{\eta} - 8\tilde{C}t} \tilde{C}^2, \quad (3.17)$$

where  $\tilde{C} = \mu^2 L_g^2 (d+6)^3 + \frac{8(d+4)}{m} \sigma^2 + \frac{8(d+4)}{m} G^2$  and  $G$  is the upper bound of the Riemannian gradient  $\text{grad}f(x)$  over the manifold  $\mathcal{M}$ . To guarantee

$$\min_{k=0, \dots, N-1} \mathbb{E}_{\mathcal{U}_k, \Xi_k} \|\bar{v}_k\|_F^2 \leq \epsilon^2 / L_g^2,$$

the parameters need to be set as:  $\mu = \mathcal{O}(\epsilon/d^{3/2})$ ,  $m = \mathcal{O}(dG^2/\epsilon^2)$ ,  $N = \mathcal{O}(1/\epsilon^2)$ . Hence, the number of calls to the stochastic zeroth-order oracle is  $\mathcal{O}(d/\epsilon^4)$ .

*Proof.* Proof of Theorem 3.3 Summing up (3.16) over  $k = 0, \dots, N-1$  and using Corollary 3.1, we have:

$$\begin{aligned} p(x_0) - \mathbb{E}_{\mathcal{U}_k, \Xi_k} p(x_k) &\geq \sum_{k=0}^{N-1} \left[ \frac{\eta_k}{2t} - \tilde{C} \right] \mathbb{E}_{\mathcal{U}_k} \|v_k\|_F^2 \geq \left[ \frac{\hat{\eta}}{4t} - 2\tilde{C} \right] \sum_{k=0}^{N-1} 2 \mathbb{E}_{\mathcal{U}_k, \Xi_k} \|v_k\|_F^2 \\ &\geq \left[ \frac{\hat{\eta}}{4t} - 2\tilde{C} \right] \sum_{k=0}^{N-1} \left[ \mathbb{E}_{\mathcal{U}_k, \Xi_k} \|\bar{v}_k\|_F^2 - t^2 \left( \mu^2 L_g^2 (d+6)^3 + \frac{8(d+4)}{m} \sigma^2 \right. \right. \\ &\quad \left. \left. + \frac{8(d+4)}{m} \|\text{grad}f(x_k)\|^2 \right) \right] \\ &\geq \left[ \frac{\hat{\eta}}{4t} - 2\tilde{C} \right] \sum_{k=0}^{N-1} \mathbb{E}_{\mathcal{U}_k, \Xi_k} \|\bar{v}_k\|_F^2 - \frac{\hat{\eta}Nt}{4} \left( \mu^2 L_g^2 (d+6)^3 + \frac{8(d+4)}{m} \sigma^2 + \frac{8(d+4)}{m} G^2 \right) \\ &\quad + 2t^2 \left( \mu^2 L_g^2 (d+6)^3 + \frac{8(d+4)}{m} \sigma^2 + \frac{8(d+4)}{m} G^2 \right)^2, \end{aligned}$$

which immediately implies the desired result (3.17).  $\square$

**Remark 3.2.** *The subproblem Eq. (3.12) is the main computational effort in Algorithm 3. Fortunately, this subproblem can be efficiently solved by a regularized semi-smooth Newton's method when  $\mathcal{M}$  takes certain forms. We refer the reader to [XLWZ18, CMMCSZ20] for more details.*

---

**Algorithm 4** Zeroth-Order Riemannian Stochastic Cubic Regularized Newton's Method (ZO-RSCRN)

---

- 1: **Input:** Initial point  $x_0$  on  $\mathcal{M}$ , smoothing parameter  $\mu$ , multi-sample parameter  $m$  and  $b$ , cubic regularization parameter  $\alpha$ , number of iteration  $N$ .
  - 2: **for**  $k = 0$  **to**  $N - 1$  **do**
  - 3:   Compute  $\bar{g}_{\mu,\xi}(x_k)$  and  $\bar{H}_{\mu,\xi}(x_k)$  based on (3.6) and (2.19) respectively.
  - 4:   Solve  $\eta_k = \operatorname{argmin}_{\eta} \hat{m}_{x_k,\alpha}(\eta)$ , where  $\hat{m}_{x,\alpha}(\eta)$  is defined in (3.18).
  - 5:   Update  $x_{k+1} = R_{x_k}(P_x(\eta_k))$ .
  - 6: **end for**
- 

### 3.4 Escaping saddle points: Zeroth-order stochastic cubic regularized Newton's method over Riemannian manifolds

In this section, we consider the problem of escaping saddle-points and converging to local minimizers in a stochastic zeroth-order Riemannian setting. Towards that, we leverage the Hessian estimator methodology developed in Section 2.3 and analyze a zeroth-order Riemannian stochastic cubic regularized Newton's method (ZO-RSCRN) for solving (3.5), which provably escapes the saddle points. Our approach is motivated by [ZZ18], where the authors proposed the minimization of function  $m_{x,\sigma}(\eta) = f(x) + \langle \operatorname{grad} f(x), \eta \rangle + \frac{1}{2} \langle P_x \circ \operatorname{Hess} f(x) \circ P_x[\eta], \eta \rangle + \frac{\alpha}{6} \|\eta\|^3$  at each iteration. The zeroth-order counterpart replaces the Riemannian gradient and Hessian with the corresponding zeroth-order estimators. The proposed ZO-RSCRN algorithm is described in Algorithm 4. In ZO-RSCRN, the function in the cubic regularized subproblem is

$$\hat{m}_{x,\alpha}(\eta) = f(x) + \langle \bar{g}_{\mu,\xi}(x), \eta \rangle + \frac{1}{2} \langle \bar{H}_{\mu,\xi}(x)[\eta], \eta \rangle + \frac{\alpha}{6} \|\eta\|^3. \quad (3.18)$$

Note that if  $\hat{\eta} = \operatorname{argmin}_{\eta} \hat{m}_{x,\alpha}(\eta)$ , then the projection  $P_x(\hat{\eta})$  is also a minimizer, because  $\bar{g}_{\mu,\xi}(x)$  and  $\bar{H}_{\mu,\xi}(x)$  only take effect on the component that is in  $T_x\mathcal{M}$ .

**Theorem 3.4.** *For manifold  $\mathcal{M}$  and function  $f : \mathcal{M} \rightarrow \mathbb{R}$  under Assumptions 2.1, 2.2 and 2.3, define  $k_{\min} := \operatorname{argmin}_k \mathbb{E}_{\mathcal{U}_k, \Xi_k} \|\eta_k\|$ , then the update in Algorithm 4 with  $\alpha \geq L_H$  satisfies:*

$$\mathbb{E} \|g_{k_{\min}+1}\| \leq \mathcal{O}(\epsilon), \text{ and } \mathbb{E} [\lambda_{\min}(\operatorname{Hess} f_{k_{\min}+1})] \geq -\mathcal{O}(\sqrt{\epsilon}), \quad (3.19)$$

given that the parameters satisfy:

$$N = \mathcal{O}\left(1/\epsilon^{3/2}\right), \quad \mu = \mathcal{O}\left(\min\left\{\frac{\epsilon}{d^{3/2}}, \sqrt{\frac{\epsilon}{d^5}}\right\}\right), \quad m = \mathcal{O}(d/\epsilon^2), \quad b = \mathcal{O}(d^4/\epsilon), \quad (3.20)$$

where  $\lambda_{\min}$  denotes the smallest eigenvalue. Hence, the zeroth-order oracle complexity is  $\mathcal{O}(d/\epsilon^{7/2} + d^4/\epsilon^{5/2})$ .

*Proof.* Proof of Theorem 3.4 Denote  $f_k = f(x_k)$ ,  $g_k = \operatorname{grad} f(x_k)$  and  $\mathbb{E} = \mathbb{E}_{\mathcal{U}_k, \Xi_k}$  for ease of notation. We first provide the global optimality conditions of subproblem Eq. (3.18) following [NP06]:

$$(\bar{H}_{\mu,\xi}(x) + \lambda^* I)\eta + \bar{g}_{\mu,\xi}(x) = 0, \quad \lambda^* = \frac{\alpha}{2} \|\eta\|, \quad \bar{H}_{\mu,\xi}(x) + \lambda^* I \succeq 0. \quad (3.21)$$

Since the parallel transport  $P_\eta$  is an isometry, we have

$$\begin{aligned}
& \|g_{k+1}\| = \|P_{\eta_k}^{-1}g_{k+1}\| \\
& = \|(P_{\eta_k}^{-1}g_{k+1} - g_k - \text{Hess}f_k[\eta_k]) + (g_k - \bar{g}_{\mu,\xi}(x_k)) \\
& \quad + (\text{Hess}f_k[\eta_k] - \bar{H}_{\mu,\xi}(x_k)[\eta_k]) + (\bar{g}_{\mu,\xi}(x_k) + \bar{H}_{\mu,\xi}(x_k)[\eta_k])\| \\
& \leq \|P_{\eta_k}^{-1}g_{k+1} - g_k - \text{Hess}f_k[\eta_k]\| + \|g_k - \bar{g}_{\mu,\xi}(x_k)\| \\
& \quad + \|\text{Hess}f_k[\eta_k] - \bar{H}_{\mu,\xi}(x_k)[\eta_k]\| + \|\bar{g}_{\mu,\xi}(x_k) + \bar{H}_{\mu,\xi}(x_k)[\eta_k]\| \\
Eq. (2.11) & \leq \frac{L_H}{2}\|\eta_k\|^2 + \|g_k - \bar{g}_{\mu,\xi}(x_k)\| \\
& \quad + \|\text{Hess}f_k[\eta_k] - \bar{H}_{\mu,\xi}(x_k)[\eta_k]\| + \|\bar{g}_{\mu,\xi}(x_k) + \bar{H}_{\mu,\xi}(x_k)[\eta_k]\| \\
Eq. (3.21) & = \frac{L_H}{2}\|\eta_k\|^2 + \|g_k - \bar{g}_{\mu,\xi}(x_k)\| + \|\text{Hess}f_k[\eta_k] - \bar{H}_{\mu,\xi}(x_k)[\eta_k]\| + \lambda^*\|\eta_k\| \\
Eq. (3.21) & \leq \frac{L_H}{2}\|\eta_k\|^2 + \|g_k - \bar{g}_{\mu,\xi}(x_k)\| + \|\text{Hess}f_k[\eta_k] - \bar{H}_{\mu,\xi}(x_k)\|_{\text{op}}\|\eta_k\| + \frac{\alpha}{2}\|\eta_k\|^2 \\
& \leq \frac{L_H}{2}\|\eta_k\|^2 + \|g_k - \bar{g}_{\mu,\xi}(x_k)\| + \frac{1}{2}\|\text{Hess}f_k - \bar{H}_{\mu,\xi}(x_k)\|_{\text{op}}^2 + \frac{1}{2}\|\eta_k\|^2 + \frac{\alpha}{2}\|\eta_k\|^2.
\end{aligned}$$

Taking expectation on both sides of the above inequality gives (by Eq. (3.8) and Eq. (2.20))

$$\mathbb{E}\|g_{k+1}\| - \sqrt{\delta_g} - \delta_H \leq \frac{1}{2}(L_H + \alpha + 1 + 2L_2\|g_k\|)\mathbb{E}\|\eta_k\|^2, \quad (3.22)$$

where  $\delta_g = \mu^2 L_g^2 (d+6)^3 + \frac{8(d+4)}{m}(G^2 + \sigma^2)$ ,  $G$  is the upper bound of  $\|\text{grad}f\|$  over  $\mathcal{M}$ , and  $\delta_H = \frac{(d+16)^4}{b}L_g + \frac{\mu^2 L_H^2}{18}(d+6)^5$ . Since  $P_{\eta_k}^{-1}$  is an isometry, we have:

$$\begin{aligned}
& \lambda_{\min}(\text{Hess}f_{k+1}) = \lambda_{\min}(P_{\eta_k}^{-1} \circ \text{Hess}f_{k+1} \circ P_{\eta_k}) \\
& \geq \lambda_{\min}(P_{\eta_k}^{-1} \circ \text{Hess}f_{k+1} \circ P_{\eta_k} - \text{Hess}f_k) \\
& \quad + \lambda_{\min}(\text{Hess}f_k - \bar{H}_{\mu,\xi}(x_k)) + \lambda_{\min}(\bar{H}_{\mu,\xi}(x_k)) \\
Eq. (2.10) & \geq -L_H\|\eta_k\| + \lambda_{\min}(\text{Hess}f_k - \bar{H}_{\mu,\xi}(x_k)) + \lambda_{\min}(\bar{H}_{\mu,\xi}(x_k)) \\
& = \lambda_{\min}(\text{Hess}f_k - \bar{H}_{\mu,\xi}(x_k)) + \lambda_{\min}(\bar{H}_{\mu,\xi}(x_k) - L_H\|\eta_k\|I) \\
Eq. (3.21) & \geq \lambda_{\min}(\text{Hess}f_k - \bar{H}_{\mu,\xi}(x_k)) - \frac{\alpha + 2L_H}{2}\|\eta_k\|.
\end{aligned}$$

Taking expectation, we obtain (by Eq. (2.20))

$$\frac{\alpha + 2L_H}{2}\mathbb{E}\|\eta_k\| \geq -(\sqrt{\delta_H} + \mathbb{E}\lambda_{\min}(\text{Hess}f_{k+1})). \quad (3.23)$$

Now we will upper bound  $\mathbb{E}\|\eta_k\|$ . From Assumption 2.2, we have

$$\begin{aligned}
\hat{f}_{x_k}(\eta_k) & \leq f(x_k) + g_k^\top \eta_k + \frac{1}{2}\eta_k^\top H_k \eta_k + \frac{L_H}{6}\|\eta_k\|^3 \\
& = \left( f(x_k) + \bar{g}_\mu(x_k)^\top \eta_k + \frac{1}{2}\eta_k^\top \bar{H}_\mu(x_k) \eta_k + \frac{L_H}{6}\|\eta_k\|^3 \right) \\
& \quad + \left( (g_k - \bar{g}_\mu(x_k))^\top \eta_k + \frac{1}{2}\eta_k^\top (H_k - \bar{H}_\mu(x_k)) \eta_k \right).
\end{aligned} \quad (3.24)$$

Using Eq. (3.21) we have

$$\begin{aligned}
& f(x_k) + \bar{g}_\mu(x_k)^\top \eta_k + \frac{1}{2} \eta_k^\top \bar{H}_\mu(x_k) \eta_k + \frac{L_H}{6} \|\eta_k\|^3 \\
&= f(x_k) - \frac{1}{2} \eta_k^\top \bar{H}_\mu(x_k) \eta_k + \left( \frac{L_H}{6} - \frac{\alpha}{2} \right) \|\eta_k\|^3 \\
&= f(x_k) - \frac{1}{2} \eta_k^\top (\bar{H}_\mu(x_k) + \frac{\alpha}{2} \|\eta_k\| I) \eta_k - \left( \frac{\alpha}{4} - \frac{L_H}{6} \right) \|\eta_k\|^3 \\
&\leq f(x_k) - \left( \frac{\alpha}{4} - \frac{L_H}{6} \right) \|\eta_k\|^3 \leq f(x_k) - \frac{\alpha}{12} \|\eta_k\|^3,
\end{aligned} \tag{3.25}$$

where the last inequality is due to  $\alpha \geq L_H$ . Moreover, by Cauchy-Schwarz inequality and Young's inequality, we have

$$\begin{aligned}
& \mathbb{E} \left[ (g_k - \bar{g}_\mu(x_k))^\top \eta_k + \frac{1}{2} \eta_k^\top (H_k - \bar{H}_\mu(x_k)) \eta_k \right] \\
&\leq \mathbb{E} \|g_k - \bar{g}_\mu(x_k)\| \|\eta_k\| + \frac{1}{2} \mathbb{E} \|H_k - \bar{H}_\mu(x_k)\|_{\text{op}} \|\eta_k\|^2 \\
&\leq \frac{32}{3\alpha} \mathbb{E} \|g_k - \bar{g}_\mu(x_k)\|^{3/2} + \frac{12}{\alpha} \mathbb{E} \|H_k - \bar{H}_\mu(x_k)\|_{\text{op}}^3 + \frac{\alpha}{24} \mathbb{E} \|\eta_k\|^3.
\end{aligned} \tag{3.26}$$

Plugging (3.25) and (3.26) to Eq. (3.24), we have

$$\mathbb{E} f_{k+1} \leq f_k - \frac{\alpha}{24} \mathbb{E} \|\eta_k\|^3 + \frac{32}{3L_H} \delta_g^{3/4} + \frac{12}{L_H} \tilde{\delta}_H, \tag{3.27}$$

where  $\tilde{\delta}_H = \tilde{C} \frac{(d+16)^6}{b^{3/2}} L_g^{1.5} + \frac{1}{27} \mu^3 L_H^3 (d+6)^{7.5}$ . Taking the sum for (3.27) over  $k = 0, \dots, N-1$ , we have

$$\frac{1}{N} \sum_{k=0}^N \mathbb{E} \|\eta_k\|^3 \leq \frac{24}{L_H} \left( \frac{f_0 - f^*}{N} + \frac{32}{3L_H} \delta_g^{3/4} + \frac{12}{L_H} \tilde{\delta}_H \right),$$

which together with (3.20) yields

$$\mathbb{E} \|\eta_{k_{\min}}\|^3 \leq \mathcal{O}(\epsilon^{3/2}), \text{ and } \mathbb{E} \|\eta_{k_{\min}}\|^2 \leq \mathcal{O}(\epsilon). \tag{3.28}$$

Combining Eq. (3.28), Eq. (3.22) and Eq. (3.23) yields (3.19).  $\square$

**Remark 3.3.** *To solve the subproblem, we implement the same Krylov subspace method as in [ABBC20], where the Riemannian Hessian and vector multiplication is approximated by Lanczos iterations. Note also that in our setting, we only require vector-vector multiplications due to the structure of our Hessian estimator in Eq. (2.14). For the purpose of brevity, we refer to [CD18, ABBC20] for a comprehensive study of this method.*

## 4 Numerical Experiments and Applications

We now explore the performance of the proposed algorithms on various simulation experiments. Finally, we demonstrate the applicability of stochastic zeroth-order Riemannian optimization for the problems of zeroth-order attacks on deep neural networks and controlling stiffness matrix in robotics. We conducted our experiments on a desktop with Intel Core 9600K CPU and NVIDIA GeForce RTX 2070 GPU.

DIMENSION	$\epsilon$	STEP SIZE	NO. ITER. ZO-RGD	AVER. NO. ITER. RGD
$15 \times 5$	$10^{-3}$	$10^{-2}$	$460 \pm 137$	442
$25 \times 15$	$10^{-3}$	$10^{-2}$	$892 \pm 99$	852
$50 \times 20$	$10^{-2}$	$5 \times 10^{-3}$	$255 \pm 26$	236

Table 2: Comparison of ZO-RGD and RGD on the Procrustes problem.

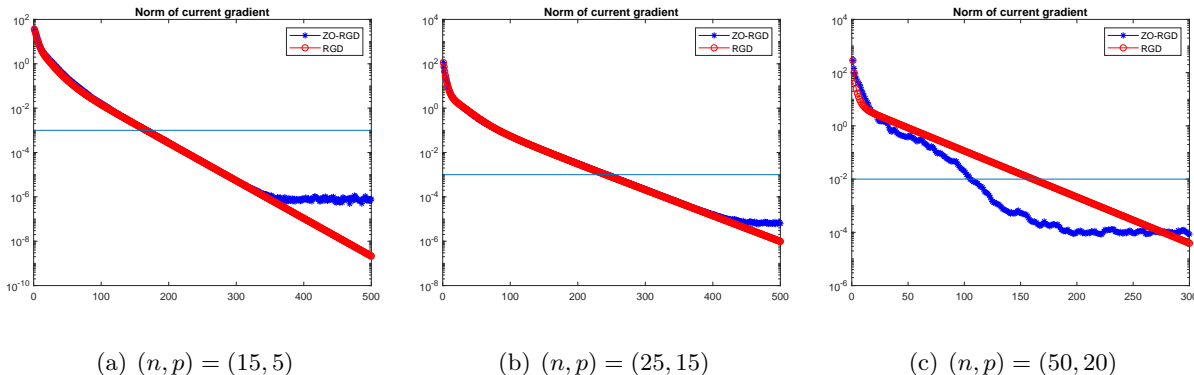


Figure 1: The convergence curve of ZO-RGD v.s. RGD. x-axis is the number of iterations and y-axis is the norm of Riemannian gradient at corresponding points. Note that our zeroth-order algorithm does not use gradient information in updates, while the graph still shows the norm of gradient to show the effectiveness of our method. The horizontal lines are the prescribed precisions.

#### 4.1 Simulation Experiments

For all the simulation experiments listed below, we plot the average result over 100 runs.

**Experiment 1: Procrustes problem [AMS09].** This is a matrix linear regression problem on a given manifold:  $\min_{X \in \mathcal{M}} \|AX - B\|_F^2$ , where  $X \in \mathbb{R}^{n \times p}$ ,  $A \in \mathbb{R}^{l \times n}$  and  $B \in \mathbb{R}^{l \times p}$ . The manifold we use is the Stiefel manifold  $\mathcal{M} = \text{St}(n, p)$ . In our experiment, we pick up different dimension  $n \times p$  and record the time cost to achieve prescribed precision  $\epsilon$ . The entries of matrix  $A$  are generated by standard Gaussian distribution. We compare our ZO-RGD (Algorithm 1) with the first-order Riemannian gradient method (RGD) on this problem. The results are shown in Table 2. For each run, we sample  $m = n \times p$  Gaussian samples for each iteration. The multi-sample version of ZO-RGD closely resembles the convergence rate of RGD, as shown in Fig. 1. These results indicate our zeroth-order method ZO-RGD is comparable with its first-order counterpart RGD, though the former one only uses zeroth-order information.

**Experiment 2: k-PCA [ZRS16, TFBJ18, ZYYF19].** k-PCA on Grassmann manifold is a Rayleigh quotient minimization problem. Given a symmetric positive definite matrix  $H \in \mathbb{R}^{n \times n}$ , we need to solve  $\min_{X \in \text{Gr}(n, p)} -\frac{1}{2} \text{Tr}(X^\top H X)$ . The Grassmann manifold  $\text{Gr}(n, p)$  is the set of  $p$ -dimensional subspaces in  $\mathbb{R}^n$ . We refer the reader to [AMS09] for more details about the Grassmann quotient manifold. This problem can be written as a finite sum problem:  $\min_{X \in \text{Gr}(n, p)} \sum_{i=1}^n -\frac{1}{2} \text{Tr}(X^\top h_i h_i^\top X)$ , where  $h_i \in \mathbb{R}^n$  and  $H = \sum_{i=1}^n h_i h_i^\top$ . We compare our ZO-RSGD algorithm (Algorithm 2) and its first-order counterpart RSGD on this problem. The results are shown in Fig. 2 (a) and (d). In our experiment, we set  $n = 100$ ,  $p = 50$ , and the matrix  $H$  is generated by  $H = AA^\top$ , where  $A \in \mathbb{R}^{n \times p}$  is a normalized randomly generated data matrix. From Fig. 2 (a) and (d), we see that the performance of ZO-RSGD is similar to its first-order counterpart



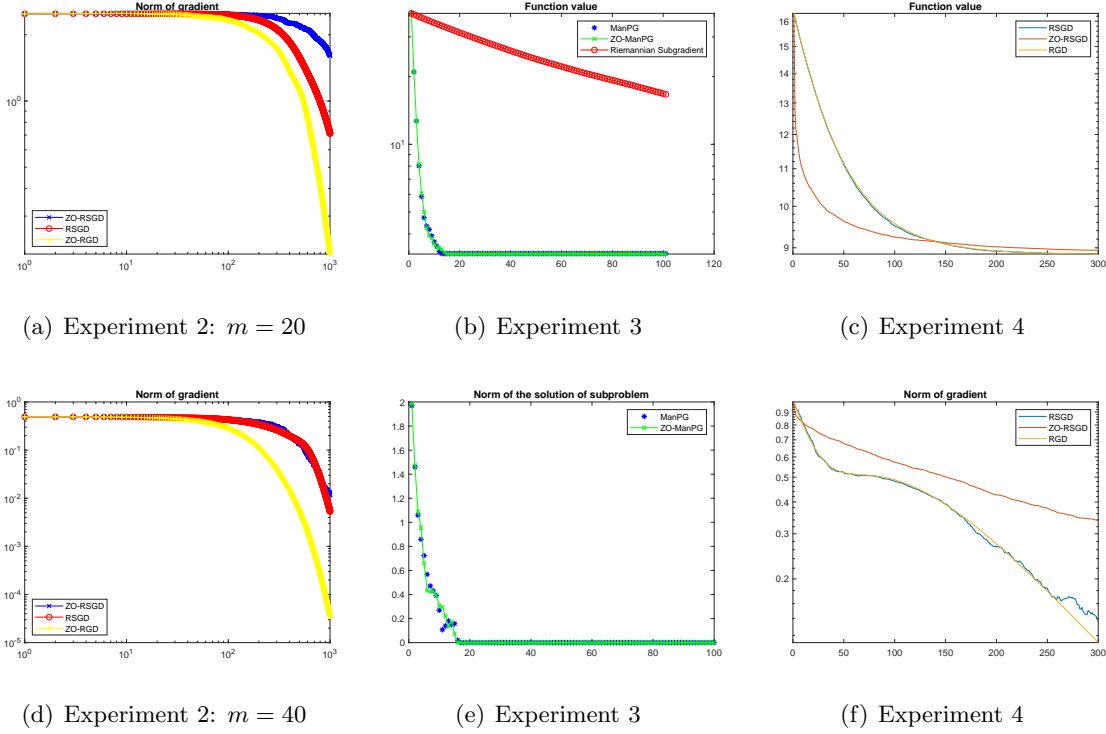


Figure 2: The convergence of three numerical experiments. The  $x$ -axis always denotes the number of iterations. Figures (a) and (d) are results for k-PCA (Experiment 2). Here three algorithms are compared: ZO-RSGD (Algorithm 2), RSGD, and ZO-RGD (Algorithm 1). Figures (b) and (e) are results for sparse PCA (Experiment 3) in which the  $y$ -axis of Figure (e) denotes the norm of  $v_k$  in (3.11) (for ManPG) and (3.12) (for ZO-ManPG), which actually measures the optimality of the problem. Here three algorithms are compared: ZO-ManPG (Algorithm 3), ManPG and Riemannian subgradient method. Figures (c) and (f) are results for Karcher mean of PSD matrices problem (Experiment 4). Here three algorithms are compared: RSGD, ZO-RSGD (Algorithm 2), and RGD.

RSGD.

**Experiment 3: Sparse PCA [JNU03, ZHT06, ZX18].** The sparse PCA problem, arising in statistics, is a Riemannian optimization problem over the Stiefel manifold with nonsmooth objective:  $\min_{X \in \text{St}(n,p)} -\frac{1}{2} \text{Tr}(X^\top A^\top A X) + \lambda \|X\|_1$ . Here,  $A \in \mathbb{R}^{m \times n}$  is the normalized data matrix. We compare our ZO-ManPG (Algorithm 3) with ManPG [CMMCSZ20] and Riemannian subgradient method [LCD<sup>+</sup>19]. In our numerical experiments, we chose  $(m, n, p) = (50, 100, 10)$ , and entries of  $A$  are drawn from Gaussian distribution and rows of  $A$  are then normalized. The comparison results are shown in Fig. 2 (b) and (e). These results show that our ZO-ManPG is comparable to its first-order counterpart ManPG and they are both much better than the Riemannian subgradient method.

**Experiment 4: Karcher mean of given PSD matrices [BI13, ZS16, KSM18].** Given a set of positive semidefinite (PSD) matrices  $\{A_i\}_{i=1}^n$  where  $A_i \in \mathbb{R}^{d \times d}$  and  $A_i \succeq 0$ , we want to calculate their Karcher mean:  $\min_{X \in \mathcal{S}_{++}^d} \frac{1}{2n} \sum_{i=1}^n (\text{dist}(X, A_i))^2$ , where  $\text{dist}(X, Y) = \|\log_m(X^{-1/2} Y X^{-1/2})\|_F$  ( $\log_m$  stands for matrix logarithm) represents the distance along the corresponding geodesic between the two points  $X, Y \in \mathcal{S}_{++}^d$ . This experiment serves as an example of optimizing geodesically convex

functions over Hadamard manifolds, with ZO-RSGD (Algorithm 2). In our numerical experiment, we take  $d = 3$  and  $n = 500$ . We compare our ZO-RSGD algorithm with its first-order counterpart RSGD and RGD. The results are shown in Fig. 2 (c) and (f), and from these results we see that ZO-RSGD is comparable to its first-order counterpart RSGD in terms of function value, though it is inferior to RSGD and RGD in terms of the size of the gradient.

**Experiment 5: Procrustes problem with ZO-RSCRN.** Here, we consider the Procrustes problem in Experiment 1 and use the ZO-RSCRN with both estimated gradients and Hessians. Following [ABBC20], we use the gradient norm as a performance measure (although the algorithm converges to local-minimizers). We use the Lanczos method (specifically Algorithm 2 from [ABBC20]) for solving the sub-problem in Step 4. Furthermore, as we are estimating the second order information, we set  $n = 6$  and  $p = 4$  and consider  $\epsilon = 10^{-3}$ . In Figure 3, (a), we plot the gradient norm versus iterations for Riemannian Stochastic Cubic-Regularized Newton method in the zeroth order and second-order setting. We notice that the zeroth-order method compares favourably to the second-order counterpart in terms of iteration complexity. Admittedly, scaling up the ZO-RSCRN method to work in higher-dimensions, based on variance reduction techniques, is an interesting problem that we plan to tackle as future work.

## 4.2 Real world applications

**Black-box stiffness control for robotics.** We now study the first motivating example discussed in Section 1.2.1 on the control of robotics with the policy parameter being the stiffness matrix  $\mathbf{K}^{\mathcal{P}} \in \mathcal{S}_{++}^d$ , see [JRBC20] for more engineering details. Mathematically, given the current position of robot  $\hat{\mathbf{p}}$  and current speed  $\dot{\mathbf{p}}$ , the task is to minimize

$$f(\mathbf{K}^{\mathcal{P}}) = w_p \|\hat{\mathbf{p}} - \mathbf{p}\|^2 + w_d \det(\mathbf{K}^{\mathcal{P}}) + w_c \text{cond}(\mathbf{K}^{\mathcal{P}}) \quad (4.1)$$

with  $\mathbf{p}$  being the new position, and  $\text{cond}$  is the condition number. With a constant external force  $\mathbf{f}^e$  applied to the system, we have the following identity which solves  $\mathbf{p}$  by  $\mathbf{K}^{\mathcal{P}}$ :  $\mathbf{f}^e = \mathbf{K}^{\mathcal{P}}(\hat{\mathbf{p}} - \mathbf{p}) - \mathbf{K}^{\mathcal{D}}\dot{\mathbf{p}}$ , where the damping matrix  $\mathbf{K}^{\mathcal{D}} = \mathbf{K}^{\mathcal{P}}$  for critical damped case. As the stiffness matrix is a positive definite matrix, the above optimization problem is a Riemannian optimization problem over the positive definite manifold (where the manifold structure is the same as the Karcher mean problem). The function  $f$  is not known analytically and following [JRBC20], we use a simulated setting for a robot (7-DOF Franka Emika Panda robot) to evaluate the function  $f$  for a given value of  $\mathbf{K}^{\mathcal{P}}$ , with the same parameters as in [JRBC20]. We compare our ZO-RGD method with Euclidean Zeroth-order gradient descent (ZO-GD) method [BG19]. We test the cases when  $d = 2$  and  $d = 3$  for minimizing function  $f$  w.r.t  $\mathbf{K}^{\mathcal{P}}$ , and the results are shown in Figure 3, (b) and (c). In our experiments, the stepsize of ZO-GD is  $3 \times 10^{-4}$  and ZO-RGD is  $10^{-3}$ . Note that for ZO-GD method, one has to project the matrix back to the positive definite set, whereas the ZO-RGD method intrinsically guarantees that the iterates are positive definite, thus is much more stable. Also, due to the fact that ZO-RGD is more stable, the stepsize of ZO-RGD can be larger than ZO-GD, which results in faster convergence.

**Zeroth-order black-box attack on Deep Neural Networks (DNNs).** We now return to the motivating example described in Section 1.2.2 and propose our black-box attack algorithm, as stated in Algorithm 6 (in Appendix D). For the sake of comparison, we also assume the architecture of the DNN is known and use “white-box” attacks based on first-order Riemannian optimization methods (Algorithm 5) and compare against the PGD attack [MMS<sup>+</sup>17], which does not explicitly enforce any constraints on the perturbed training data. For simplicity, we assume the manifold is a sphere. That is, we assume that the perturbation set  $S$  is given by  $S(R) = \{\delta : \|\delta\|_2 = R\}$ , where  $R$  is the radius of the sphere. This is consistent with the optimal  $\ell_2$ -norm attack studied in the literature [LHL15]. Furthermore, the sphere constraint guarantees that the perturbed image is always in a

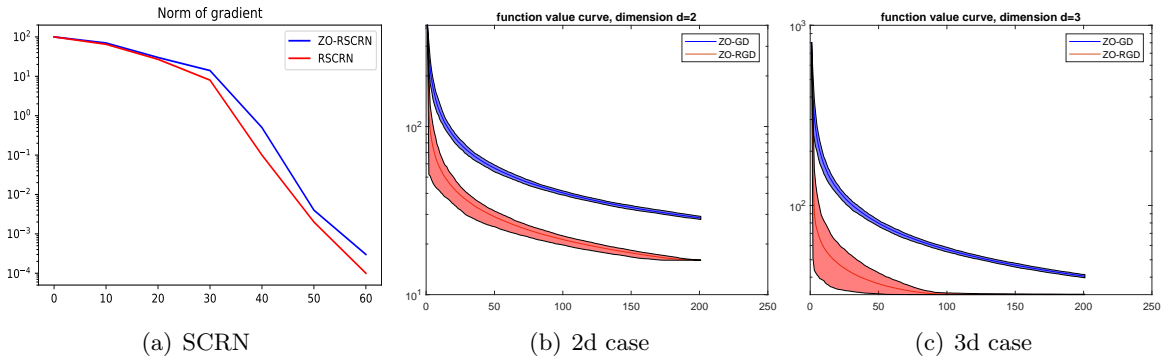


Figure 3: Figure (a) corresponds to Experiment 5. Figures (b) and (c) correspond to the experiments on the robotic minimization function in (4.1). The  $x$ -axis in all figures correspond to iteration number.

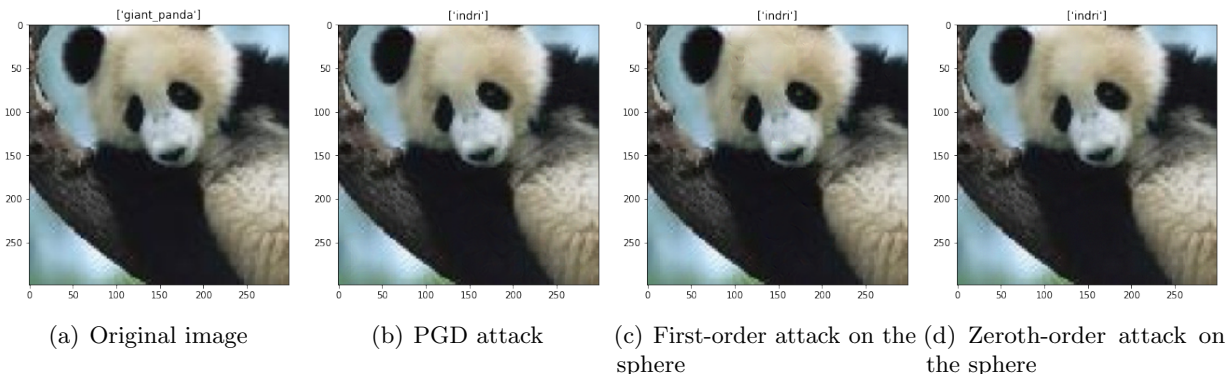


Figure 4: The attack on giant panda picture [DDS<sup>+</sup>09]. (a): the original image; (b): the PGD attack with a small diameter; (c) Riemannian attack (Algorithm 5) on the sphere with the same diameter; (d): Riemannian zeroth-order attack (Algorithm 6). 'Indri' refers to the class to which the original image is misclassified to.

certain distance from the original image. We start our zeroth-order attack from a perturbation and maximize the loss function on the sphere. For the black-box method, to accelerate the convergence, we use Euclidean zeroth-order optimization to find an appropriate initial perturbation (Algorithm 7). It is worth noting that the zeroth-order attack in [CZS<sup>+</sup>17, TTC<sup>+</sup>19] has a non-smooth objective function, which has  $\mathcal{O}(n^3/\epsilon^3)$  complexity to guarantee convergence [NS17], whereas the complexity needed for our method is  $\mathcal{O}(d/\epsilon^2)$ .

We first tested our method on the giant panda picture in the Imagenet data set [DDS<sup>+</sup>09], with the network structure the Inception v3 network structure [SVI<sup>+</sup>16]. The attack radius in our algorithm is proportional to the  $\ell_2$  norm of the original image. Both white-box and black-box Riemannian attacks are successful, which means that they both converge to images that lie in a different image class (i.e. with a different label), see Figure 4. We also tested Algorithms 5 and 6 on the CIFAR10 dataset, and the network structure we used is the VGG net [SZ14]. The corresponding results are provided in Appendix D.

## 5 Conclusions

In this paper, we proposed zeroth-order algorithms for solving Riemannian optimization over submanifolds embedded in Euclidean space in which only noisy function evaluations are available for the objective. These algorithms adopt new estimators of the Riemannian gradient and Hessian from noisy objective function evaluations, based on a Riemannian version of the Gaussian smoothing technique. The proposed estimators overcome the difficulty of the non-linearity of the manifold constraint and the issues that arise in using Euclidean Gaussian smoothing techniques when the function is defined only over the manifold. The iteration complexity and oracle complexity of the proposed algorithms are analyzed for obtaining an appropriately defined  $\epsilon$ -stationary point or  $\epsilon$ -approximate local minimum. The established complexities are independent of the dimension of the ambient Euclidean space and only depend on the intrinsic dimension of the manifold. Numerical experiments demonstrated that the proposed zeroth-order algorithms are comparable to their first-order counterparts.

## A Geodesically Convex Problem

In this section we consider the smooth problem (3.5) where  $f$  is geodesically convex. The definition of geodesic convexity is given below (see, e.g., [ZS16]).

**Definition A.1.** *A function  $f : \mathcal{M} \rightarrow \mathbb{R}$  is geodesically convex if for all  $x, y \in \mathcal{M}$ , there exists a geodesic  $\gamma$  such that  $\gamma(0) = x$ ,  $\gamma(1) = y$  and  $\forall t \in [0, 1]$  we have  $f(\gamma(t)) \leq (1-t)f(x) + tf(y)$ .*

It can be shown that this definition is equivalent to,  $f(\text{Exp}_x(\eta)) \geq f(x) + \langle g_x, \eta \rangle_x$ ,  $\forall \eta \in T_x \mathcal{M}$ , where  $g_x$  is a subgradient of  $f$  at  $x$ ,  $\text{Exp}$  is the exponential mapping, and  $\langle \cdot, \cdot \rangle_x$  is the inner product in  $T_x \mathcal{M}$  induced by Riemannian metric  $d(\cdot, \cdot)$ . When  $f$  is smooth, we have  $g_x = \text{grad}f(x)$ , the Riemannian gradient at  $x$ . It is known that geodesically convex function is a constant on compact manifolds. Therefore, in this subsection, we assume that  $\mathcal{M}$  is an Hadamard manifold [BO69, Gro78], and  $\mathcal{X}$  is a bounded and geodesically convex subset of  $\mathcal{M}$ .

**Assumption A.1.** *The subset  $\mathcal{X}$  of Hadamard manifold  $\mathcal{M}$  is bounded by diameter  $D$ , and the sectional curvature is lower bounded by  $\varrho$ . The function  $F(x, \xi)$  is geodesically convex w.r.t.  $x \in \mathcal{M}$ , almost everywhere for  $\xi$  (and hence  $f$  is geodesically convex).*

The following lemma from [ZS16] is useful for our subsequent analysis. Here  $\mathcal{P}_{\mathcal{X}}$  denotes the projection onto  $\mathcal{X}$ , i.e.,  $\mathcal{P}_{\mathcal{X}}(x) := \{y \in \mathcal{X} : d(x, y) = \inf_{z \in \mathcal{X}} d(x, z)\}$ .

**Lemma A.1** ([ZS16]). *For any Riemannian manifold  $\mathcal{M}$  where the sectional curvature is lower bounded by  $\varrho$  and any points  $x, x_s \in \mathcal{M}$ , the update*

$$x_{s+1} = \mathcal{P}_{\mathcal{X}}(\text{Exp}_{x_s}(-\eta_s g_s))$$

*satisfies:  $\langle -g_s, x - x_s \rangle \leq \frac{1}{2\eta_s}(d^2(x_s, x) - d^2(x_{s+1}, x)) + \frac{\zeta(\varrho, d(x_s, x))\eta_s}{2}\|g_s\|^2$ , where  $d(\cdot, \cdot)$  is the Riemannian metric defined globally on  $\mathcal{M}$ , and  $\zeta(\varrho, c) := c\sqrt{|\varrho|}/\tanh(c\sqrt{|\varrho|})$ .*

In this subsection, we consider the following algorithm, which is a special case of Algorithm 2.

$$x_{k+1} = \mathcal{P}_{\mathcal{X}}(\text{Exp}_{x_k}(-\eta_k \bar{g}_{\mu, \xi}(x_k))). \quad (\text{A.1})$$

We now present our result for obtaining an  $\epsilon$ -optimal solution of (3.5).

**Theorem A.1.** *Let the manifold  $\mathcal{M}$  and the function  $f : \mathcal{M} \rightarrow \mathbb{R}$  satisfy Assumptions 2.1, 2.3, and A.1. Suppose Algorithm 2 is run with the update in Eq. (A.1) and with  $\eta_k = 1/L_g$ . Denote  $\Delta_k = \mathbb{E}_{u_k, \Xi_k}(f(x_k) - f^*)$ . To have  $\min_{1 \leq k \leq t} \Delta_k \leq \epsilon$ , we need the smoothing parameter  $\mu$ , number of sampling at each iteration  $m$  and the number of iteration  $t$  to be respectively of order:*

$$\mu = \mathcal{O}(\sqrt{\epsilon}/d^{3/2}), \quad m = \mathcal{O}(d/\epsilon), \quad t = \mathcal{O}(1/\epsilon). \quad (\text{A.2})$$

Hence, the zeroth-order oracle complexity is  $N = mt = \mathcal{O}(d/\epsilon^2)$ .

*Proof.* Proof of Theorem A.1 From Assumption 2.1 we have that:

$$f(x_{k+1}) - f(x_k) \leq -\eta_k \langle \text{grad}f(x_k), \bar{g}_{\mu, \xi}(x_k) \rangle + \frac{L_g}{2} \eta_k^2 \|\bar{g}_{\mu, \xi}(x_k)\|^2.$$

Taking  $\eta_k = \frac{1}{L_g}$ , we have

$$\begin{aligned} f(x_{k+1}) - f(x_k) &\leq \frac{1}{2L_g} \left( -2 \langle \text{grad}f(x_k), \bar{g}_{\mu, \xi}(x_k) \rangle + \|\bar{g}_{\mu, \xi}(x_k)\|^2 \right) \\ &= \frac{1}{2L_g} \left( \|\bar{g}_{\mu, \xi}(x_k) - \text{grad}f(x_k)\|^2 - \|\text{grad}f(x_k)\|^2 \right). \end{aligned}$$

Taking expectation with respect to  $u_k$  on both sides of the inequality above and taking  $m \geq 16(d+4)$ , we have (by Eq. (3.8))

$$\begin{aligned} &\mathbb{E}_{u_k} f(x_{k+1}) - f(x_k) \\ &\leq \frac{1}{2L_g} \left( \mu^2 L_g^2 (d+6)^3 + \frac{8(d+4)}{m} \sigma^2 + \left( \frac{8(d+4)}{m} - 1 \right) \|\text{grad}f(x_k)\|^2 \right) \\ &\leq \frac{\mu^2 L_g (d+6)^3}{2} + \frac{4(d+4)}{m L_g} \sigma^2 - \frac{1}{4L_g} \|\text{grad}f(x_k)\|^2. \end{aligned} \quad (\text{A.3})$$

Now considering the geodesic convexity and Lemma A.1, we have

$$f(x_{k+1}) - f^* \leq \langle -\bar{g}_{\mu, \xi}(x_k), \text{Exp}_{x_k}^{-1}(x^*) \rangle \leq \frac{L_g}{2} (d^2(x_k, x^*) - d^2(x_{k+1}, x^*)) + \frac{\zeta(\varrho, D) \|\bar{g}_{\mu, \xi}(x_k)\|^2}{2L_g}. \quad (\text{A.4})$$

From Lemma 3.1 we have

$$\begin{aligned} \mathbb{E} \|\bar{g}_{\mu, \xi}(x_k)\|^2 &\leq 2\mathbb{E} \|\bar{g}_{\mu, \xi}(x_k) - \text{grad}f(x_k)\|^2 + 2\mathbb{E} \|\text{grad}f(x_k)\|^2 \\ &\leq 2\mu^2 L_g^2 (d+6)^3 + \frac{16(d+4)}{m} \sigma^2 + \left( \frac{16(d+4)}{m} + 2 \right) \|\text{grad}f(x_k)\|^2. \end{aligned} \quad (\text{A.5})$$

Now take the expectation w.r.t.  $u_k$  for both sides of (A.4), and combine with (A.5), we have

$$\Delta_{k+1} \leq \frac{L_g}{2} (d^2(x_k, x^*) - d^2(x_{k+1}, x^*)) + \frac{\zeta(\varrho, D)}{2L_g} \left( 2\mu^2 L_g^2 (d+6)^3 + \frac{16(d+4)}{m} \sigma^2 + 3\|\text{grad}f(x_k)\|^2 \right). \quad (\text{A.6})$$

Multiplying (A.6) with  $\frac{1}{6\zeta(\varrho, D)}$ , and sum up with Eq. (A.3), we have

$$\left( 1 + \frac{1}{6\zeta} \right) \Delta_{k+1} - \Delta_k \leq \frac{L_g}{12\zeta} (d^2(x_k, x^*) - d^2(x_{k+1}, x^*)) + \mu^2 L_g (d+6)^3 + \frac{16(d+4)}{3m L_g} \sigma^2.$$

Summing it over  $k = 0, \dots, t-1$  we have

$$\Delta_t - \Delta_0 + \frac{1}{6\zeta} \sum_{k=1}^t \Delta_k \leq \frac{L_g}{12\zeta} d^2(x_0, x^*) + (\mu^2 L_g (d+6)^3 + \frac{16(d+4)}{3mL_g} \sigma^2) t.$$

Equivalently, we have

$$\frac{1}{t} \sum_{k=1}^t \Delta_k \leq \frac{L_g}{2t} d^2(x_0, x^*) + 6\zeta(\mu^2 L_g (d+6)^3 + \frac{16(d+4)}{3mL_g} \sigma^2) + \frac{6\zeta}{t} \Delta_0,$$

which together with (A.2) yields  $\min_{1 \leq k \leq t} \Delta_k \leq \epsilon$ .  $\square$

## B Proof of Remark 3.1

*Proof.* Proof of the improved bound Eq. (3.3)

Since  $\bar{g}_\mu(x) = \frac{1}{m} \sum_{i=1}^m g_{\mu,i}(x)$ , we have (denote  $\mathcal{U} = \{u_1, \dots, u_m\}$ ):

$$\begin{aligned} & \mathbb{E}_{\mathcal{U}} \|\bar{g}_\mu(x) - \text{grad}f(x)\|^2 \\ & \leq 2\mathbb{E}_{\mathcal{U}} \|\bar{g}_\mu(x) - \mathbb{E}_{\mathcal{U}} \bar{g}_\mu(x)\|^2 + 2\|\mathbb{E}_{\mathcal{U}} \bar{g}_\mu(x) - \text{grad}f(x)\|^2 \\ & = 2\mathbb{E}_{\mathcal{U}} \left\| \frac{1}{m} \sum_{i=1}^m [g_{\mu,i}(x) - \mathbb{E}_{\mathcal{U}} g_{\mu,i}(x)] \right\|^2 + 2 \left\| \frac{1}{m} \sum_{i=1}^m [\mathbb{E}_{\mathcal{U}} g_{\mu,i}(x) - \text{grad}f(x)] \right\|^2 \\ & = \frac{2}{m^2} \mathbb{E}_{\mathcal{U}} \sum_{i=1}^m \|g_{\mu,i}(x) - \mathbb{E}_{\mathcal{U}} g_{\mu,i}(x)\|^2 + \frac{2}{m^2} \left\| \sum_{i=1}^m [\mathbb{E}_{\mathcal{U}} g_{\mu,i}(x) - \text{grad}f(x)] \right\|^2 \\ & \leq \frac{2}{m} \mathbb{E}_{u_1} \|g_{\mu,1}(x) - \mathbb{E}_{\mathcal{U}} g_{\mu,1}(x)\|^2 + 2\|\mathbb{E}_{u_1} g_{\mu,1}(x) - \text{grad}f(x)\|^2 \\ & \leq \frac{2}{m} \mathbb{E}_{u_1} \|g_{\mu,1}(x)\|^2 + \frac{\mu^2 L_g^2}{2} (d+3)^3 \\ & \leq \frac{\mu^2}{m} L_g^2 (d+6)^3 + \frac{4(d+4)}{m} \|\text{grad}f(x)\|^2 + \frac{\mu^2 L_g^2}{2} (d+3)^3 \\ & \leq \mu^2 L_g^2 (d+6)^3 + \frac{4(d+4)}{m} \|\text{grad}f(x)\|^2, \end{aligned}$$

where the second equality is from the fact that  $u_i$  and  $u_j$  are independent when  $i \neq j$ .  $\square$

*Proof.* Proof of (3.4) Following the  $L_g$ -retraction-smooth, we have:  $f(x_{k+1}) \leq f(x_k) - \eta_k \langle \bar{g}_\mu(x), \text{grad}f(x_k) \rangle + \frac{\eta_k^2 L_g}{2} \|\bar{g}_\mu(x)\|^2$ . Taking  $\eta_k = \hat{\eta} = 1/L_g$ , we have

$$\begin{aligned} f(x_{k+1}) & \leq f(x_k) - \eta_k \langle \bar{g}_\mu(x), \text{grad}f(x_k) \rangle + \frac{\eta_k^2 L_g}{2} \|\bar{g}_\mu(x)\|^2 \\ & = f(x_k) + \frac{1}{2L_g} (\|\bar{g}_\mu(x) - \text{grad}f(x)\|^2 - \|\text{grad}f(x)\|^2). \end{aligned}$$

Now take the expectation for the random variables at the iteration  $k$  on both sides, we have

$$\begin{aligned} \mathbb{E}_k f(x_{k+1}) & \leq f(x_k) + \frac{1}{2L_g} (\mathbb{E}_k \|\bar{g}_\mu(x) - \text{grad}f(x)\|^2 - \|\text{grad}f(x)\|^2) \\ \text{Eq. (3.8)} & \leq f(x_k) + \frac{1}{2L_g} \left( \mu^2 L_g^2 (d+6)^3 + \left( \frac{4(d+4)}{m} - 1 \right) \|\text{grad}f(x)\|^2 \right). \end{aligned}$$

By choosing  $m \geq 8(d+4)$ , summing the above inequality over  $k = 0, \dots, N$  gives (3.4).  $\square$

## C Proof of Lemma 3.1

*Proof.* For the sake of notation, here we denote  $\mathbb{E} = \mathbb{E}_{u_0}$ . From (3.6) we have

$$\mathbb{E}(\|g_{\mu,\xi}(x)\|^2) = \frac{1}{\mu^2} \mathbb{E} [(F(R_x(\mu u, \xi)) - F(x, \xi))^2 \|u\|^2]. \quad (\text{C.1})$$

From Assumption 2.1 we have

$$\begin{aligned} & (F(R_x(\mu u, \xi)) - F(x, \xi))^2 \\ &= (F(R_x(\mu u, \xi)) - F(x, \xi) - \mu \langle \text{grad}F(x, \xi), u \rangle + \mu \langle \text{grad}F(x, \xi), u \rangle)^2 \\ &\leq 2 \left( \frac{L_g}{2} \mu^2 \|u\|^2 \right)^2 + 2\mu^2 \langle \text{grad}F(x, \xi), u \rangle^2. \end{aligned} \quad (\text{C.2})$$

Combining (C.1) and (C.2) yields

$$\begin{aligned} \mathbb{E}(\|g_{\mu,\xi}(x)\|^2) &\leq \frac{\mu^2}{2} L_g^2 \mathbb{E}(\|u\|^6) + 2\mathbb{E}(\|\langle \text{grad}F(x, \xi), u \rangle u\|^2) \\ (\text{Corollary 2.1}) &\leq \frac{\mu^2}{2} L_g^2 (d+6)^3 + 2\mathbb{E}(\|\langle \text{grad}F(x, \xi), u \rangle u\|^2). \end{aligned} \quad (\text{C.3})$$

Denote our  $d$ -dimensional tangent space as  $\mathcal{X}$ . Without loss of generality, suppose  $\mathcal{X}$  is the subspace generated by projecting onto the first  $d$  coordinates, i.e.,  $\forall x \in \mathcal{X}$ , the last  $(n-d)$  elements of  $x$  are zeros. Also for brevity, denote  $g = \text{grad}F(x, \xi)$ . Use  $x_i$  to denote the  $i$ -th coordinate of  $u_0$ , and  $\kappa(d)$  denote the normalization constant for  $d$ -dimensional Gaussian distribution. For simplicity, denote  $x = (x_1, \dots, x_d)$ . We have

$$\begin{aligned} \mathbb{E}(\|\langle \text{grad}F(x, \xi), u \rangle u\|^2) &= \frac{1}{\kappa} \int_{\mathbb{R}^n} \langle \text{grad}F(x, \xi), u \rangle^2 \|u\|^2 e^{-\frac{1}{2}\|u_0\|^2} du_0 \\ &= \frac{1}{\kappa(d)} \int_{\mathbb{R}^d} \left( \sum_{i=1}^d g_i x_i \right)^2 \left( \sum_{i=1}^d x_i^2 \right) e^{-\frac{1}{2} \sum_{i=1}^d x_i^2} dx_1 \cdots dx_d, \\ &= \frac{1}{\kappa(d)} \int_{\mathbb{R}^d} \langle g, x \rangle^2 \|x\|^2 e^{-\frac{1}{2}\|x\|^2} dx = \frac{1}{\kappa(d)} \int_{\mathbb{R}^d} \|x\|^2 e^{-\frac{\tau}{2}\|x\|^2} \langle g, x \rangle^2 e^{-\frac{1-\tau}{2}\|x\|^2} dx \\ &\leq \frac{2}{\kappa(d)\tau e} \int_{\mathbb{R}^d} \langle g, x \rangle^2 e^{-\frac{1-\tau}{2}\|x\|^2} dx = \frac{2}{\kappa(d)\tau(1-\tau)^{1+d/2} e} \int_{\mathbb{R}^d} \langle g, x \rangle^2 e^{-\frac{1}{2}\|x\|^2} dx \\ &= \frac{2}{\tau(1-\tau)^{1+d/2} e} \|g\|^2 \leq (d+4) \|g\|^2, \end{aligned} \quad (\text{C.4})$$

where the last  $n-d$  dimensions of  $u_0$  are integrated to be one, the first inequality is due to the following fact:  $x^p e^{-\frac{\tau}{2}x^2} \leq (\frac{2}{\tau e})^{p/2}$ , and the second inequality follows by setting  $\tau = \frac{2}{(d+4)}$ . From Assumption 2.3, we have

$$\mathbb{E}_\xi \|\text{grad}F(x, \xi)\|^2 \leq 2\mathbb{E}_\xi \|\text{grad}F(x, \xi) - \text{grad}f(x)\|^2 + 2\|\text{grad}f(x)\|^2 \leq 2\sigma^2 + 2\|\text{grad}f(x)\|^2. \quad (\text{C.5})$$

Combining (C.3), (C.4) and (C.5) yields

$$\begin{aligned} \mathbb{E}_\xi [\mathbb{E}_{u_0}(\|g_{\mu,\xi}(x)\|^2)] &\leq \mathbb{E}_\xi \left[ \frac{\mu^2}{2} L_g^2 (d+6)^3 + 2(d+4) \|\text{grad}F(x, \xi)\|^2 \right] \\ &\leq \frac{\mu^2}{2} L_g^2 (d+6)^3 + 4(d+4)(\sigma^2 + \|\text{grad}f(x)\|^2). \end{aligned} \quad (\text{C.6})$$

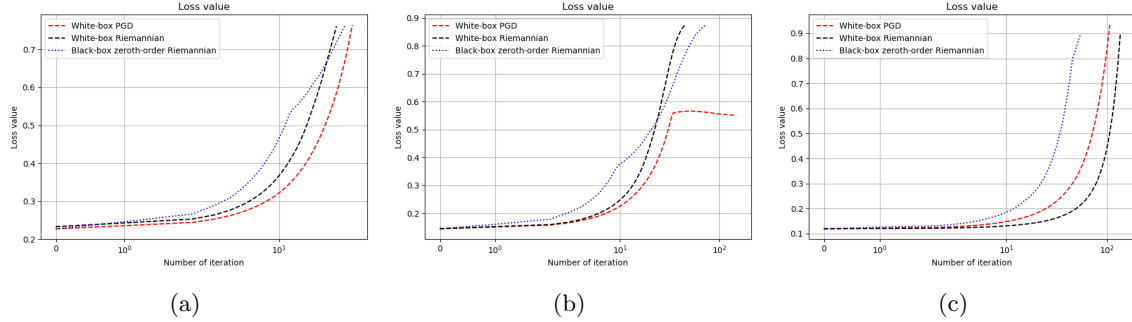


Figure 5: Loss function versus the iteration numbers. We observe that the loss function increases while performing our attacks. The three figures correspond to the last three rows in Figure 6. For the failed the PGD attack, we notice that the function value stuck in the middle and then decreased, while white and black-box Riemannian attack increased loss value successfully.

Finally, we have

$$\begin{aligned}
& \mathbb{E}_{\mathcal{U}, \xi} \|\bar{g}_{\mu, \xi}(x) - \text{grad}f(x)\|^2 \\
& \leq 2\mathbb{E}_{\mathcal{U}, \xi} \|\bar{g}_{\mu, \xi}(x) - \mathbb{E}_{\mathcal{U}, \xi} g_{\mu}(x)\|^2 + 2\|\mathbb{E}_{\mathcal{U}, \xi} g_{\mu}(x) - \text{grad}f(x)\|^2 \\
& \leq \frac{2}{m} \mathbb{E}_{u_1, \xi_1} \|g_{\mu_1, \xi_1}(x) - \mathbb{E}_{u_1, \xi_1} g_{\mu_1, \xi_1}(x)\|^2 + \frac{\mu^2 L_g^2}{2} (d+3)^3 \\
& \leq \frac{2}{m} \mathbb{E}_{u_1, \xi_1} \|g_{\mu_1, \xi_1}(x)\|^2 + \frac{\mu^2 L_g^2}{2} (d+3)^3 \\
& \leq \frac{2}{m} \left( \frac{\mu^2 L_g^2}{2} (d+6)^3 + 4(d+4)[\|\text{grad}f(x)\|^2 + \sigma^2] \right) + \frac{\mu^2 L_g^2}{2} (d+3)^3 \\
& \leq \mu^2 L_g^2 (d+6)^3 + \frac{8(d+4)}{m} \sigma^2 + \frac{8(d+4)}{m} \|\text{grad}f(x)\|^2,
\end{aligned}$$

where the second inequality is from Proposition 2.1 and the fourth inequality is from (C.6).  $\square$

## D Implementation Details of Black-Box Attacks

Here we provide our white and black-box Riemannian attack algorithm in Algorithms 5 and 6, respectively. For the black-box attack, to accelerate convergence, we introduce a pre-attack step to search for a sufficiently large loss value on the prescribed sphere, still in a black-box manner (only use the function value). For further acceleration of the black-box attack, the hierarchical attack [CZS<sup>+</sup>17] and the auto-encoder technique [TTC<sup>+</sup>19] might be applicable. The attack results on CIFAR-10 images are shown in Figure 6. We also provide the loss function curve in Figure 5. Again the network structure we used is the VGG net [SZ14]. It can be seen from Figure 6 that our black-box attack yields similar attack result as the PGD attack, however PGD failed in one of the five images, and the loss function curve in Figure 5 indicates that the loss function curve of PGD attack stagnated. This may be due to an inappropriate choice of parameters for the PGD attack. Nevertheless this experiment shows the ability of our Riemannian attack method. We also mention here that the Riemannian attack is in general slower than the PGD attack due to the multi-sampling technique discussed in Remark 3.1.



---

**Algorithm 5** White-box attack via Riemannian optimization

---

- 1: **Input:** Original image  $\tilde{x}$ , original label  $y_0$ , radius of the attack region  $R$ , step size  $\eta_k$ , some convergence criterion.
  - 2: Randomly sample  $\delta$  s.t.  $\|\delta\| = R$ .
  - 3: Set the initial point  $x_0 = \tilde{x} + \delta$ ,  $k = 0$ .
  - 4: **repeat**
  - 5:   Update  $x_{k+1} = R_{x_k}(\eta_k \text{grad}_x L(\theta, x, y))$ ,  $k = k + 1$ .
  - 6: **until** Convergence criterion is met.
- 

---

**Algorithm 6** Black-box attack via Riemannian zeroth-order optimization

---

- 1: **Input:** Original image  $\tilde{x}$ , original label  $y_0$ , radius of the attack region  $R$ , step size  $\eta_k$ , smoothing parameter  $\mu$ , number of multi-sample  $m$ , some convergence criterion.
  - 2: Obtain  $\delta$  (initial perturbation) by pre-training steps (Algorithm 7).
  - 3: Set the initial point  $x_0 = \tilde{x} + \delta$ ,  $k = 0$ .
  - 4: **repeat**
  - 5:   Sample  $m$  standard Gaussian random matrix  $u_i$  on  $T_{x_k}\mathcal{M}$ .
  - 6:   Set the random oracle  $\bar{g}_\mu(x_k)$  by (3.6).
  - 7:   Update  $x_{k+1} = R_{x_k}(\eta_k \bar{g}_\mu(x_k))$ ,  $k = k + 1$ .
  - 8: **until** Convergence criterion is met.
- 

---

**Algorithm 7** Pre-training step for black-box attack

---

- 1: **Input:** Original image  $\tilde{x}$ , original label  $y_0$ , radius of the attack region  $R$ , step size  $\eta_k$ , smoothing parameter  $\mu$ , number of multi-sample  $m$ .
  - 2:  $x_0 = \tilde{x}$ .
  - 3: **repeat**
  - 4:   Sample  $m$  standard Gaussian random matrix  $u_i$ .
  - 5:   Set the random oracle  $\bar{g}_\mu(x_k)$  by (3.6), with  $g_{\mu_i}(x) = \frac{f(x+\mu u_i) - f(x)}{\mu} u_i$
  - 6:   Update  $x_{k+1} = x_k + \eta_k \bar{g}_\mu(x_k)$ .
  - 7:   Update  $\delta = x_{k+1} - \tilde{x}$ ,  $k = k + 1$ .
  - 8: **until**  $\|\delta\| \geq R$
  - 9:  $\delta = \frac{\delta}{\|\delta\|} R$
- 

## Acknowledgments.

JL and SM acknowledge the support by NSF grants DMS-1953210 and CCF-2007797. KB and SM acknowledge the support by UC Davis CeDAR (Center for Data Science and Artificial Intelligence Research) Innovative Data Science Seed Funding Program. JL also wants to thank Tesi Xiao for helpful discussions on black-box attacks.

## References

- [ABBC20] Naman Agarwal, Nicolas Boumal, Brian Bullins, and Coralia Cartis, *Adaptive regularization with cubics on manifolds*, Mathematical Programming (to appear) (2020).
- [AH17] Charles Audet and Warren Hare, *Derivative-free and blackbox optimization*, Springer, 2017.

- [AMS09] P-A Absil, Robert Mahony, and Rodolphe Sepulchre, *Optimization algorithms on matrix manifolds*, Princeton University Press, 2009.
- [BA11] Nicolas Boumal and Pierre Absil, *RTRMC: A Riemannian trust-region method for low-rank matrix completion*, Advances in neural information processing systems, 2011, pp. 406–414.
- [BAC18] Nicolas Boumal, Pierre-Antoine Absil, and Coralia Cartis, *Global rates of convergence for nonconvex optimization on manifolds*, IMA Journal of Numerical Analysis **39** (2018), no. 1, 1–33.
- [BEB17] Tamir Bendory, Yonina C Eldar, and Nicolas Boumal, *Non-convex phase retrieval from STFT measurements*, IEEE Transactions on Information Theory **64** (2017), no. 1, 467–484.
- [BFM17] Glaydston C Bento, Orizon P Ferreira, and Jefferson G Melo, *Iteration-complexity of gradient, subgradient and proximal point methods on Riemannian manifolds*, Journal of Optimization Theory and Applications **173** (2017), no. 2, 548–562.
- [BG19] Krishnakumar Balasubramanian and Saeed Ghadimi, *Zeroth-order nonconvex stochastic optimization: Handling constraints, high-dimensionality, and saddle-points*, arXiv preprint arXiv:1809.06474 (2019), 651–676.
- [BI13] Dario A Bini and Bruno Iannazzo, *Computing the Karcher mean of symmetric positive definite matrices*, Linear Algebra and its Applications **438** (2013), no. 4, 1700–1710.
- [BO69] Richard L Bishop and Barrett O’Neill, *Manifolds of negative curvature*, Transactions of the American Mathematical Society **145** (1969), 1–49.
- [Bon13] Silvere Bonnabel, *Stochastic gradient descent on Riemannian manifolds*, IEEE Transactions on Automatic Control **58** (2013), no. 9, 2217–2229.
- [BSBA14] Pierre B Borckmans, S Easter Selvan, Nicolas Boumal, and P-A Absil, *A Riemannian subgradient algorithm for economic dispatch with valve-point effect*, Journal of Computational and Applied Mathematics **255** (2014), 848–866.
- [Bur10] Christopher JC Burges, *Dimension reduction: A guided tour*, Now Publishers Inc, 2010.
- [Car92] Manfredo Perdigao do Carmo, *Riemannian geometry*, Birkhäuser, 1992.
- [CD18] Yair Carmon and John C Duchi, *Analysis of Krylov subspace solutions of regularized non-convex quadratic problems*, Advances in Neural Information Processing Systems, 2018, pp. 10705–10715.
- [CDMS20] S. Chen, Z. Deng, S. Ma, and A. M.-C. So, *Manifold proximal point algorithms for dual principal component pursuit and orthogonal dictionary learning*, <https://arxiv.org/abs/2005.02356> (2020).
- [CLC<sup>+</sup>18] Minhao Cheng, Thong Le, Pin-Yu Chen, Jinfeng Yi, Huan Zhang, and Cho-Jui Hsieh, *Query-efficient hard-label black-box attack: An optimization-based approach*, arXiv preprint arXiv:1807.04457 (2018).

- [CM17] Frédéric Chazal and Bertrand Michel, *An introduction to topological data analysis: fundamental and practical aspects for data scientists*, arXiv preprint arXiv:1710.04019 (2017).
- [CMMCSZ20] Shixiang Chen, Shiqian Ma, Anthony Man-Cho So, and Tong Zhang, *Proximal gradient method for nonsmooth optimization over the Stiefel manifold*, SIAM Journal on Optimization **30** (2020), no. 1, 210–239.
- [CMYZ20] HanQin Cai, Daniel Mckenzie, Wotao Yin, and Zhenliang Zhang, *Zeroth-order regularized optimization (zoro): Approximately sparse gradients and adaptive sampling*, arXiv preprint arXiv:2003.13001 (2020).
- [CS16] Anoop Cherian and Suvrit Sra, *Riemannian dictionary learning and sparse coding for positive definite matrices*, IEEE transactions on neural networks and learning systems **28** (2016), no. 12, 2859–2871.
- [CSA15] Amit Chattopadhyay, Suviseshamuthu Easter Selvan, and Umberto Amato, *A derivative-free Riemannian Powell’s method, minimizing hartley-entropy-based ICA contrast*, IEEE transactions on neural networks and learning systems **27** (2015), no. 9, 1983–1990.
- [CSV09] Andrew Conn, Katya Scheinberg, and Luis Vicente, *Introduction to derivative-free optimization*, vol. 8, SIAM, 2009.
- [CZS<sup>+</sup>17] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh, *ZOO: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models*, Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security, ACM, 2017, pp. 15–26.
- [DDS<sup>+</sup>09] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, *ImageNet: A Large-Scale Hierarchical Image Database*, CVPR09, 2009.
- [DET17] Danny Drieß, Peter Englert, and Marc Toussaint, *Constrained bayesian optimization of combined interaction force/task space controllers for manipulations*, 2017 IEEE International Conference on Robotics and Automation (ICRA), IEEE, 2017, pp. 902–907.
- [DHS<sup>+</sup>13] Persi Diaconis, Susan Holmes, Mehrdad Shahshahani, et al., *Sampling from a manifold*, Advances in modern statistical theory and applications: a Festschrift in honor of Morris L. Eaton, Institute of Mathematical Statistics, 2013, pp. 102–125.
- [DJWW15] John C Duchi, Michael I Jordan, Martin J Wainwright, and Andre Wibisono, *Optimal rates for zero-order convex optimization: The power of two function evaluations*, IEEE Transactions on Information Theory **61** (2015), no. 5, 2788–2806.
- [Fra18] Peter I Frazier, *A tutorial on Bayesian optimization*, arXiv preprint arXiv:1807.02811 (2018).
- [FT19] Robert Simon Fong and Peter Tino, *Stochastic derivative-free optimization on Riemannian manifolds*, arXiv preprint arXiv:1908.06783 (2019).

- [GKK<sup>+</sup>19] Daniel Golovin, John Karro, Greg Kochanski, Chansoo Lee, and Xingyou Song, *Gradientless descent: High-dimensional zeroth-order optimization*, arXiv preprint arXiv:1911.06317 (2019).
- [GL13] Saeed Ghadimi and Guanghui Lan, *Stochastic first-and zeroth-order methods for nonconvex stochastic programming*, SIAM Journal on Optimization **23** (2013), no. 4, 2341–2368.
- [Gro78] Mikhael Gromov, *Manifolds of negative curvature*, Journal of Differential Geometry **13** (1978), no. 2, 223–230.
- [GSS14] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy, *Explaining and harnessing adversarial examples*, arXiv preprint arXiv:1412.6572 (2014).
- [HSH17] 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** (2017), no. 1, 48–62.
- [Hsu02] Elton P Hsu, *Stochastic analysis on manifolds*, vol. 38, American Mathematical Soc., 2002.
- [JNR12] Kevin G Jamieson, Robert Nowak, and Ben Recht, *Query complexity of derivative-free optimization*, Advances in Neural Information Processing Systems, 2012, pp. 2672–2680.
- [JNU03] I. Jolliffe, N. Trendafilov, and M. Uddin, *A modified principal component technique based on the LASSO*, Journal of computational and Graphical Statistics **12** (2003), no. 3, 531–547.
- [JR20] Noémie Jaquier and Leonel Rozo, *High-dimensional Bayesian optimization via nested Riemannian manifolds*, Advances in Neural Information Processing Systems **33** (2020).
- [JR20] Noémie Jaquier, Leonel Rozo, Sylvain Calinon, and Mathias Bürger, *Bayesian optimization meets Riemannian manifolds in robot learning*, Conference on Robot Learning, PMLR, 2020, pp. 233–246.
- [Kac20] Oleg Kachan, *Persistent homology-based projection pursuit*, Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2020, pp. 856–857.
- [KGB16] Artiom Kovnatsky, Klaus Glashoff, and Michael M Bronstein, *MADMM: a generic algorithm for non-smooth optimization on manifolds*, European Conference on Computer Vision, Springer, 2016, pp. 680–696.
- [KSM18] Hiroyuki Kasai, Hiroyuki Sato, and Bamdev Mishra, *Riemannian stochastic recursive gradient algorithm*, International Conference on Machine Learning, 2018, pp. 2516–2524.
- [LCD<sup>+</sup>19] Xiao Li, Shixiang Chen, Zengde Deng, Qing Qu, Zhihui Zhu, and Anthony Man Cho So, *Nonsmooth optimization over Stiefel manifold: Riemannian subgradient methods*, arXiv preprint arXiv:1911.05047 (2019).

- [LHL15] Chunchuan Lyu, Kaizhu Huang, and Hai-Ning Liang, *A unified gradient regularization family for adversarial examples*, 2015 IEEE International Conference on Data Mining, IEEE, 2015, pp. 301–309.
- [LKPS16] Chun-Liang Li, Kirthevasan Kandasamy, Barnabás Póczos, and Jeff Schneider, *High dimensional bayesian optimization via restricted projection pursuit models*, Artificial Intelligence and Statistics, 2016, pp. 884–892.
- [LLSD20] Lizhen Lin, Drew Lazar, Bayan Sarpabayeva, and David B. Dunson, *Robust optimization and inference on manifolds*, arXiv preprint arXiv:2006.06843 (2020).
- [LMW19] Jeffrey Larson, Matt Menickelly, and Stefan M Wild, *Derivative-free optimization methods*, Acta Numerica **28** (2019), 287–404.
- [LO14] Rongjie Lai and Stanley Osher, *A splitting method for orthogonality constrained problems*, Journal of Scientific Computing **58** (2014), no. 2, 431–449.
- [LOT19] Jacob Leygonie, Steve Oudot, and Ulrike Tillmann, *A framework for differential calculus on persistence barcodes*, arXiv preprint arXiv:1910.00960 (2019).
- [LSTZD17] 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** (2017), no. 519, 1261–1273.
- [LV07] John A Lee and Michel Verleysen, *Nonlinear dimensionality reduction*, Springer Science & Business Media, 2007.
- [Mat65] J Matyas, *Random optimization*, Automation and Remote control **26** (1965), no. 2, 246–253.
- [MHB<sup>+</sup>16] Alonso Marco, Philipp Hennig, Jeannette Bohg, Stefan Schaal, and Sebastian Trimpe, *Automatic LQR tuning based on Gaussian process global optimization*, 2016 IEEE international conference on robotics and automation (ICRA), IEEE, 2016, pp. 270–277.
- [MHM18] Leland McInnes, John Healy, and James Melville, *UMAP: Uniform manifold approximation and projection for dimension reduction*, arXiv preprint arXiv:1802.03426 (2018).
- [MHST17] Alonso Marco, Philipp Hennig, Stefan Schaal, and Sebastian Trimpe, *On the design of LQR kernels for efficient controller learning*, 2017 IEEE 56th Annual Conference on Decision and Control (CDC), IEEE, 2017, pp. 5193–5200.
- [MK18] Mojmir Mutny and Andreas Krause, *Efficient high dimensional Bayesian optimization with additivity and quadrature fourier features*, Advances in Neural Information Processing Systems, 2018, pp. 9005–9016.
- [MKJS19] Bamdev Mishra, Hiroyuki Kasai, Pratik Jawanpuria, and Atul Saroop, *A Riemannian gossip approach to subspace learning on Grassmann manifold*, Machine Learning (2019), 1–21.
- [MMS<sup>+</sup>17] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu, *Towards deep learning models resistant to adversarial attacks*, arXiv preprint arXiv:1706.06083 (2017).

- [Moc94] Jonas Mockus, *Application of Bayesian approach to numerical methods of global and stochastic optimization*, Journal of Global Optimization **4** (1994), no. 4, 347–365.
- [Moc12] ———, *Bayesian approach to global optimization: theory and applications*, vol. 37, Springer Science & Business Media, 2012.
- [Nes11] Yu Nesterov, *Random gradient-free minimization of convex functions*, Technical Report. Center for Operations Research and Econometrics (CORE), Catholic University of Louvain (2011).
- [NM65] John A Nelder and Roger Mead, *A simplex method for function minimization*, The computer journal **7** (1965), no. 4, 308–313.
- [NP06] Yurii Nesterov and Boris T Polyak, *Cubic regularization of Newton method and its global performance*, Mathematical Programming **108** (2006), no. 1, 177–205.
- [NS17] Yurii Nesterov and Vladimir Spokoiny, *Random gradient-free minimization of convex functions*, Foundations of Computational Mathematics **17** (2017), no. 2, 527–566.
- [NY83] Arkadi S. Nemirovski and David B. Yudin, *Problem complexity and method efficiency in optimization*, Wiley & Sons (1983).
- [OGW18] ChangYong Oh, Efstratios Gavves, and Max Welling, *BOCK: Bayesian optimization with cylindrical kernels*, International Conference on Machine Learning, 2018, pp. 3868–3877.
- [RB19] Raúl Rabadán and Andrew J Blumberg, *Topological data analysis for genomics and evolution: Topology in biology*, Cambridge University Press, 2019.
- [Rio09] Emmanuel Rio, *Moment inequalities for sums of dependent random variables under projective conditions*, Journal of Theoretical Probability **22** (2009), no. 1, 146–163.
- [RSBC18] Paul Rolland, Jonathan Scarlett, Ilija Bogunovic, and Volkan Cevher, *High-dimensional Bayesian optimization via additive models with overlapping groups*, International Conference on Artificial Intelligence and Statistics, 2018, pp. 298–307.
- [Spa05] James C Spall, *Introduction to stochastic search and optimization: Estimation, simulation, and control*, vol. 65, John Wiley & Sons, 2005.
- [SQW16] Ju Sun, Qing Qu, and John Wright, *Complete dictionary recovery over the sphere I: Overview and the geometric picture*, IEEE Transactions on Information Theory **63** (2016), no. 2, 853–884.
- [SQW18] ———, *A geometric analysis of phase retrieval*, Foundations of Computational Mathematics **18** (2018), no. 5, 1131–1198.
- [SSW<sup>+</sup>15] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas, *Taking the human out of the loop: A review of Bayesian optimization*, Proceedings of the IEEE **104** (2015), no. 1, 148–175.
- [Ste72] Charles Stein, *A bound for the error in the Normal approximation to the distribution of a sum of dependent random variables*, Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability, Volume 2: Probability Theory, The Regents of the University of California, 1972.

- [SVI<sup>+</sup>16] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna, *Rethinking the inception architecture for computer vision*, Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 2818–2826.
- [SZ14] Karen Simonyan and Andrew Zisserman, *Very deep convolutional networks for large-scale image recognition*, arXiv preprint arXiv:1409.1556 (2014).
- [SZS<sup>+</sup>13] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus, *Intriguing properties of neural networks*, arXiv preprint arXiv:1312.6199 (2013).
- [TFBJ18] Nilesh Tripuraneni, Nicolas Flammarion, Francis Bach, and Michael I Jordan, *Averaging stochastic gradient descent on Riemannian manifolds*, Conference On Learning Theory, 2018, pp. 650–687.
- [TTC<sup>+</sup>19] Chun-Chen Tu, Paishun Ting, Pin-Yu Chen, Sijia Liu, Huan Zhang, Jinfeng Yi, Cho-Jui Hsieh, and Shin-Ming Cheng, *AutoZOOM: Autoencoder-based zeroth order optimization method for attacking black-box neural networks*, Proceedings of the AAAI Conference on Artificial Intelligence, vol. 33, 2019, pp. 742–749.
- [Van13] Bart Vandereycken, *Low-rank matrix completion by Riemannian optimization*, SIAM Journal on Optimization **23** (2013), no. 2, 1214–1236.
- [WDBS18] Yining Wang, Simon Du, Sivaraman Balakrishnan, and Aarti Singh, *Stochastic zeroth-order optimization in high dimensions*, International Conference on Artificial Intelligence and Statistics, 2018, pp. 1356–1365.
- [WFT20] Linnan Wang, Rodrigo Fonseca, and Yuandong Tian, *Learning search space partition for black-box optimization using monte carlo tree search*, Advances in Neural Information Processing Systems **33** (2020).
- [WHZ<sup>+</sup>16] Ziyu Wang, Frank Hutter, Masrour Zoghi, David Matheson, and Nando de Freitas, *Bayesian optimization in a billion dimensions via random embeddings*, Journal of Artificial Intelligence Research **55** (2016), 361–387.
- [WLC<sup>+</sup>20] Z. Wang, B. Liu, S. Chen, S. Ma, L. Xue, and H. Zhao, *A manifold proximal linear method for sparse spectral clustering with application to single-cell RNA sequencing data analysis*, <https://arxiv.org/abs/2007.09524> (2020).
- [WMX20] B. Wang, S. Ma, and L. Xue, *Riemannian stochastic proximal gradient methods for nonsmooth optimization over the Stiefel manifold*, <https://arxiv.org/pdf/2005.01209.pdf> (2020).
- [WS04] K. Q. Weinberger and L. K. Saul, *Unsupervised learning of image manifolds by semidenite programming*, CVPR, 2004.
- [WS19] Melanie Weber and Suvrit Sra, *Nonconvex stochastic optimization on manifolds via Riemannian Frank-Wolfe methods*, arXiv preprint arXiv:1910.04194 (2019).
- [XLWZ18] Xiantao Xiao, Yongfeng Li, Zaiwen Wen, and Liwei Zhang, *A regularized semi-smooth Newton method with projection steps for composite convex programs*, Journal of Scientific Computing **76** (2018), no. 1, 364–389.

- [YCL19] Kai Yuan, Iordanis Chatzinikolaïdis, and Zhibin Li, *Bayesian optimization for whole-body control of high-degree-of-freedom robots through reduction of dimensionality*, IEEE Robotics and Automation Letters **4** (2019), no. 3, 2268–2275.
- [YZS14] Wei Hong Yang, Lei-Hong Zhang, and Ruyi Song, *Optimality conditions for the nonlinear programming problems on Riemannian manifolds*, Pacific Journal of Optimization **10** (2014), no. 2, 415–434.
- [ZHT06] H. Zou, T. Hastie, and R. Tibshirani, *Sparse principal component analysis*, J. Comput. Graph. Stat. **15** (2006), no. 2, 265–286.
- [ZMZ20] J. Zhang, S. Ma, and S. Zhang, *Primal-dual optimization algorithms over Riemannian manifolds: an iteration complexity analysis*, Mathematical Programming Series A **184** (2020), 445–490.
- [ZRS16] Hongyi Zhang, Sashank J Reddi, and Suvrit Sra, *Riemannian SVRG: Fast stochastic optimization on Riemannian manifolds*, Advances in Neural Information Processing Systems, 2016, pp. 4592–4600.
- [ZS16] Hongyi Zhang and Suvrit Sra, *First-order methods for geodesically convex optimization*, Conference on Learning Theory, 2016, pp. 1617–1638.
- [ZX18] Hui Zou and Lingzhou Xue, *A selective overview of sparse principal component analysis*, Proceedings of the IEEE **106** (2018), no. 8, 1311–1320.
- [ZYYF19] 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).
- [ZZ18] Junyu Zhang and Shuzhong Zhang, *A cubic regularized Newton’s method over Riemannian manifolds*, arXiv preprint arXiv:1805.05565 (2018).



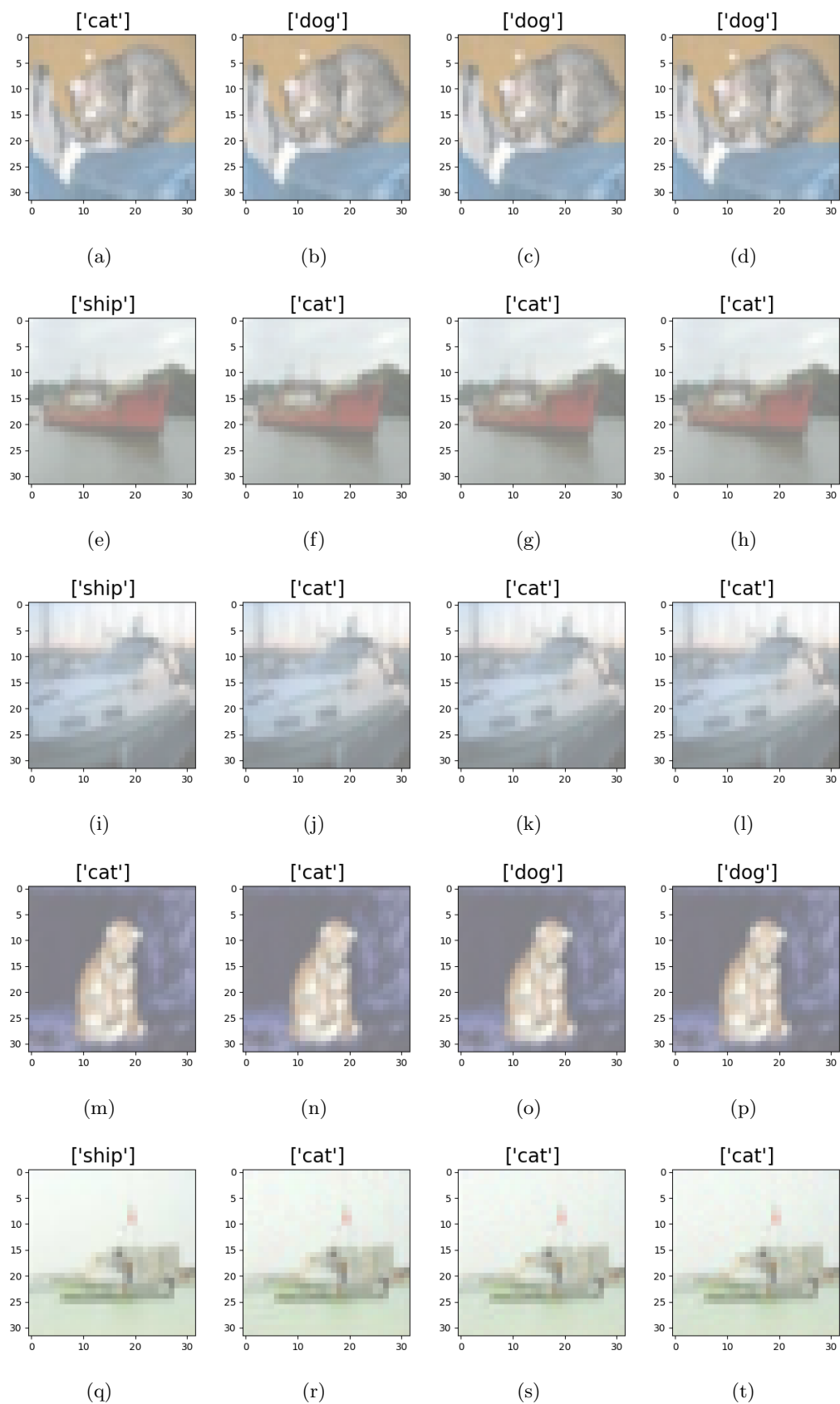


Figure 6: The attack on CIFAR10 picture. From left to right columns: the original image; the PGD attack with a small diameter; white box Riemannian attack on the sphere with the same diameter; black box Riemannian attack on the sphere with the same diameter. Notice that for the figure in the fourth row, the PGD attack failed while the Riemannian attacks succeeded. The diameter is set to be 0.01 times the norm of the original images.