

---

# On Riemannian Stochastic Approximation Schemes with Fixed Step-Size

---

**Alain Durmus**

Université Paris-Saclay  
ENS Paris-Saclay, CNRS  
Centre Borelli, F-91190 Gif-sur-Yvette, France  
alain.durmus@ens-paris-saclay.fr

**Pablo Jiménez**

CMAP, École Polytechnique  
Institut Polytechnique de Paris  
pablo.jimenez-moreno@polytechnique.edu

**Éric Moulines**

CMAP, École Polytechnique, CNRS  
Institut Polytechnique de Paris  
eric.moulines@polytechnique.edu

**Salem Said**

Laboratoire IMS, Université de Bordeaux  
salem.said@u-bordeaux.fr

## Abstract

This paper studies fixed step-size stochastic approximation (SA) schemes, including stochastic gradient schemes, in a Riemannian framework. It is motivated by several applications, where geodesics can be computed explicitly, and their use accelerates crude Euclidean methods. A fixed step-size scheme defines a family of time-homogeneous Markov chains, parametrized by the step-size. Here, using this formulation, non-asymptotic performance bounds are derived, under Lyapunov conditions. Then, for any step-size, the corresponding Markov chain is proved to admit a unique stationary distribution, and to be geometrically ergodic. This result gives rise to a family of stationary distributions indexed by the step-size, which is further shown to converge to a Dirac measure, concentrated at the solution of the problem at hand, as the step-size goes to 0. Finally, the asymptotic rate of this convergence is established, through an asymptotic expansion of the bias, and a central limit theorem.

## 1 INTRODUCTION

This paper deals with the study of fixed step-size Stochastic Approximation (SA) algorithms (Robbins

and Monro, 1951; Kushner and Yin, 2003; Polyak and Juditsky, 1992), defined on a Riemannian manifold  $\Theta$  with metric  $\mathbf{g}$ . Specifically, consider the problem

$$\begin{aligned} \text{find } \theta \in \Theta \text{ satisfying } h(\theta) = 0, \\ \text{for a vector field } h : \Theta \rightarrow T\Theta, \end{aligned} \tag{1}$$

where  $T\Theta$  denotes the tangent bundle of  $\Theta$ , and  $h$  is only accessible through an oracle returning noisy estimates. The setting where  $h = -\text{grad } f$  is of particular interest for minimizing a smooth function  $f : \Theta \rightarrow \mathbb{R}$ . In the Euclidean setting, Stochastic Gradient Descent (SGD) and its variants are now common methods for solving this problem (Bottou, 2010; Bottou and Bousquet, 2008). However, it should be stressed that (1) encompasses several other applications in stochastic optimization, reinforcement learning or maximum likelihood estimation, such as online Expectation Maximization algorithms (Cappé and Moulines, 2009), policy gradient (Baxter and Bartlett, 2001) or Q-learning (Jaakkola et al., 1993). Minimization over a Riemannian manifold or its general formulation (1) arises in many applications: Principal Component Analysis (Edelman et al., 1998), dictionary recovery (Sun et al., 2017), matrix completion (Boumal and Absil, 2011), smooth semidefinite programs (Boumal et al., 2016), tensor factorization (Ishteva et al., 2011), and Riemannian barycenter estimation (Said and Manton, 2019; Arnaudon et al., 2012). This has motivated the development of a comprehensive framework for stochastic optimization problems on Riemannian manifolds. One of the first contributions in this field is Bonnabel (2013), which derives asymptotic convergence results for SA on Riemannian manifolds. Non-asymptotic results are obtained by Zhang and Sra (2016) for a geodesically convex function  $f$ . This study has been followed and completed by Zhang et al. (2016); Sato et al. (2019)

which introduce and analyze a Riemannian counterpart of the Stochastic Variance Reduced Gradient (SVRG) algorithm. Since then, many existing methods or results from the Euclidean case have been considered in a Riemannian setting. For example, [Khuzani and Li \(2017\)](#) suggest a Riemannian stochastic primal-dual algorithm and most recently [Tripuraneni et al. \(2018\)](#) study an averaged version of Riemannian SGD.

In this paper, we are interested in the study of fixed step-size SA methods of the form

$$\theta_{n+1} = \text{proj}_{\mathcal{S}} [\text{Exp}_{\theta_n} \{\eta H_{\theta_n}(X_{n+1})\}] , \quad (2)$$

where  $H_{\theta_n}(X_{n+1}) = h(\theta_n) + e_{\theta_n}(X_{n+1})$ .

In (2),  $\eta > 0$  is a step-size,  $(X_n)_{n \in \mathbb{N}^*}$  is an  $(\mathcal{F}_n)_{n \in \mathbb{N}^*}$ -adapted process, defined on a filtered probability space, with values in a measurable space  $(\mathcal{X}, \mathcal{X})$ , and  $e : \Theta \times \mathcal{X} \rightarrow T\Theta$  is a measurable function, such that  $\theta \mapsto e_{\theta}(x)$  is a vector field over  $\Theta$ , for any  $x \in \mathcal{X}$ . In addition,  $\text{Exp}_{\theta} : T_{\theta}\Theta \rightarrow \Theta$  is the Riemannian exponential mapping and  $\text{proj}_{\mathcal{S}} : \Theta \rightarrow \mathcal{S}$  is a projection-like operator onto a subset  $\mathcal{S} \subset \Theta$ . This recursion is a natural extension of Euclidean SA, akin to the Robbins-Monroe algorithm, in a Riemannian setting.

In the Euclidean setting, the study of fixed step-size SA, and in particular SGD, has recently attracted much attention, see e.g. [Ma et al. \(2018\)](#); [Vaswani et al. \(2019\)](#); [Dieuleveut et al. \(2017\)](#); [Bach \(2020\)](#); [Bach and Moulines \(2011\)](#). Indeed, first of all, the step-size  $\eta$  is the only parameter to tune, in contrast to the case where a decreasing sequence of step-sizes is used in (2). Furthermore, the forgetting of the initial condition is exponentially fast ([Nedić and Bertsekas, 2001](#); [Needell et al., 2014](#)).

We aim to show, in a general Riemannian framework, that the use of (2) provides a good solution for (1). To this end, we establish non-asymptotic and asymptotic properties of  $(\theta_n)_{n \in \mathbb{N}}$ , in the limit  $\eta \rightarrow 0$ . Our contributions can be summarized as follows.

- (1) We derive non-asymptotic bounds, for the convergence of  $(\theta_n)_{n \in \mathbb{N}}$  to approximate solutions of (1), under general Lyapunov assumptions and mild assumptions on the manifold  $\Theta$  and the subset  $\mathcal{S}$ .
- (2) Under additional regularity conditions, we show that  $(\theta_n)_{n \in \mathbb{N}}$ , as a Markov chain, admits a unique stationary distribution  $\mu^{\eta}$  and is geometrically ergodic, *i.e.* converges to  $\mu^{\eta}$  exponentially fast.
- (3) We study the limiting behavior of the family  $(\mu^{\eta})_{\eta > 0}$  as  $\eta \rightarrow 0$ . In particular, we show that if (1) admits a unique solution  $\theta^*$  and other suitable conditions hold, this family converges to the Dirac measure at  $\theta^*$ . In addition, we asymptotically quantify this convergence, through a central limit theorem. Precisely, we prove that after a  $\eta^{-1/2}$ -rescaling, this family of

stationary distributions converges weakly to a normal distribution as  $\eta \rightarrow 0$ . These results illustrate the exponential forgetting of initial condition of the scheme and that, at stationarity, the iterates  $(\theta_n)_{n \in \mathbb{N}}$  stay in a  $\mathcal{O}(\eta^{1/2})$ -neighborhood of  $\theta^*$ . In addition, they can be understood as generalizations to Riemannian spaces of [Pflug \(1986, Theorem 1\)](#) and [Dieuleveut et al. \(2017, Theorem 4\)](#).

(4) We apply our results to SGD. In particular, we establish the first non-asymptotic convergence bounds for strongly geodesically convex functions, without boundedness assumptions on the manifold  $\Theta$ .

(5) Finally, we introduce and prove the convergence of an SGD scheme to compute the Riemannian barycenter, also known as the Karcher mean, of distributions on Hadamard manifolds. To the authors' knowledge, our contribution on this topic is one of the few without boundedness assumptions on the distribution.

For ease of reading, all the assumptions are gathered in the supplement Section S1. In the derivation of our results, we use crucially the fact that  $(\theta_n)_{n \in \mathbb{N}}$  defines a Markov chain in  $\Theta$ , under mild conditions. This interpretation has been successfully used in several papers dealing with the convergence of SA or SGD in Euclidean spaces; see e.g. [Benveniste et al. \(1990\)](#); [Kushner and Huang \(1981\)](#); [Fort and Pagès \(1999\)](#); [Pflug \(1986\)](#).

We consider a more general setting and milder conditions in comparison with most other studies in the field. Indeed, most papers do not consider the general SA framework, but only the case  $h = -\text{grad} f$ , dealing with SGD and its variants. To the authors' knowledge, only [Bonnabel \(2013\)](#); [Durmus et al. \(2020\)](#) tackle the general SA problem (1). Our main contribution, compared to these two works, is to deal with the fixed step-size setting. Besides, our study considers general geodesically complete Riemannian manifolds which encompass Hadamard spaces, which have been the primary focus for [Zhang and Sra \(2016\)](#); [Zhang et al. \(2016\)](#); [Tripuraneni et al. \(2018\)](#).

Furthermore, a majority of the previous studies on SGD in a Riemannian space (see e.g. [Zhang and Sra \(2016\)](#); [Zhang et al. \(2016\)](#); [Tripuraneni et al. \(2018\)](#); [Alimisis et al. \(2020\)](#); [Han and Gao \(2020\)](#)), are purely local in nature, because of the assumption that  $(\theta_n)_{n \in \mathbb{N}}$  stays almost surely in a (fixed and deterministic) compact and geodesically convex subset of  $\Theta$ . For example, note that all the convergence results derived in [Zhang and Sra \(2016\)](#) depend on the diameter of the compact in which  $(\theta_n)_{n \in \mathbb{N}}$  is assumed to stay. This assumption rarely holds in practice, and is quite difficult to verify in theory. It strongly limits the applicability of many results in the literature over the past few years. On the contrary, our results do not suffer from this problem,

and can all be applied either on a compact or non-compact Riemannian manifold. As a result, we consider a new SA method to estimate the Karcher mean of a distribution  $\pi$  on  $\Theta$ , see [Arnaudon et al. \(2012\)](#); [Le \(2004\)](#); [Zhang and Sra \(2016\)](#); [Iannazzo and Porcelli \(2018\)](#), for which we derive non-asymptotic convergence bounds without boundedness conditions on the support of  $\pi$ .

**Notations** For any  $\theta \in \Theta$  and  $v, w \in T_\theta\Theta$ , denote by  $\mathbf{g}_\theta(v, w) = \langle v, w \rangle_\theta$  and its corresponding norm by  $\mathbf{g}_\theta(v, v) = \|v\|_\theta^2$ .  $\rho_\Theta : \Theta \times \Theta \rightarrow \mathbb{R}_+$  denotes the distance associated with the Riemannian metric  $\mathbf{g}$ . For any  $\theta_0 \in \Theta, r > 0$ , set  $B(\theta_0, r) = \{\theta_1 \in \Theta : \rho_\Theta(\theta_0, \theta_1) < r\}$ , the open ball centered at  $\theta_0$  with radius  $r$ . Similarly, we define closed balls in  $\Theta$  by  $\bar{B}(\theta_0, r) = \{\theta_1 \in \Theta : \rho_\Theta(\theta_0, \theta_1) \leq r\}$ .

For a smooth function  $g : \Theta \rightarrow \mathbb{R}$ , we denote by  $\text{grad } g$  its Riemannian gradient ([Lee, 2019](#), p. 27) and by  $\text{Hess } g$  its Riemannian, or covariant, Hessian ([Lee, 2019](#), Example 4.22). For a curve  $\gamma : I \rightarrow \Theta, T_{\gamma(t_0)}^Y : T_{\gamma(t_0)}\Theta \rightarrow T_{\gamma(t_1)}\Theta$  stands for the parallel transport map associated to the Levi-Civita connection along  $\gamma$  from  $\gamma(t_0)$  to  $\gamma(t_1)$  ([Lee, 2019](#), Equation 4.22). Moreover, for any  $\theta \in \Theta$ , under the assumption that  $\Theta$  is complete, consider the Riemannian exponential map  $\text{Exp}_\theta : T_\theta\Theta \rightarrow \Theta$ , see [Lee \(2019, Proposition 5.19\)](#). This map projects a vector from the tangent space  $T_\theta\Theta$  onto the manifold  $\Theta$ , following a geodesic curve.

## 2 CONSTANT STEPSIZE ANALYSIS FOR A CONSTRAINED SCHEME

### 2.1 Main Results

In this section, we study the Stochastic Approximation scheme (2), which is constrained on a subset  $S \subset \Theta$ . The following assumption on the manifold  $\Theta$  and  $S$  is considered all along this paper and allows us to rigorously define  $\text{proj}_S$ .

**A1.** Assume one of the following conditions.

(i)  $\Theta$  is a Hadamard manifold, i.e. a complete, simply connected Riemannian manifold with non-positive sectional curvature. In addition,  $S$  is a closed geodesically convex subset of  $\Theta$  with non-empty interior.

(ii)  $\Theta$  is a complete, connected Riemannian manifold and  $S = \Theta$ .

Note that under **A1**, the exponential map  $\text{Exp} : T\Theta \rightarrow \Theta$  is well-defined, see [Lee \(2019, Theorem 6.19\)](#). Under **A1-(i)**, [Sturm \(2003, Proposition 2.6\)](#) shows that there exists  $\text{proj}_S : \Theta \rightarrow S$  which is the Riemannian counterpart of the Euclidean projection onto a closed convex subset. More precisely,  $\text{proj}_S$  is the unique mapping from  $\Theta$  to  $S$  such that for any  $\theta \in \Theta$ ,

$\rho_\Theta(\text{proj}_S(\theta), \theta) = \inf_{\theta' \in S} \rho_\Theta(\theta', \theta)$ . Under **A1-(ii)**, we simply set  $\text{proj}_S = \text{Id}$ .

Recall that the recursion (2) only uses a noisy estimate  $H_\theta$  of the mean field  $h(\theta)$ , for any  $\theta \in \Theta$ . We assume the following conditions on the noise to ensure convergence.

**MD 1.** The sequence  $(X_n)_{n \in \mathbb{N}^*}$  is independent and identically distributed (i.i.d.). In addition, for any  $\theta \in \Theta$ ,  $\mathbb{E}[e_\theta(X_1)] = 0$  and there exist  $\sigma_0^2, \sigma_1^2 > 0$  such that for any  $\theta \in S$ ,  $\mathbb{E}[\|e_\theta(X_1)\|_\theta^2] \leq \sigma_0^2 + \sigma_1^2 \|h(\theta)\|_\theta^2$ .

**MD1** is referred to as the martingale difference setting which implies that  $(\theta_n)_{n \in \mathbb{N}}$  is a time-homogeneous  $(\mathcal{F}_n)_{n \in \mathbb{N}}$ -Markov chain, for which we denote by  $Q_\eta$  its corresponding Markov kernel. The Euclidean counterpart of this assumption consists in replacing the Riemannian norm with the Euclidean one, which is usual in Stochastic Approximation (see [Duflo, 1997](#)).

**MD2.** (i)  $\mathbb{P}$ -almost surely, the vector field  $\theta \mapsto e_\theta(X_1)$  is continuous on  $\Theta$ .

(ii) For any  $\theta \in \Theta$ ,  $\text{Leb}_\theta$  and the distribution of  $e_\theta(X_1)$  are mutually absolutely continuous, where  $\text{Leb}_\theta$  stands for the Lebesgue measure on  $T_\theta\Theta$ .

**MD2** ensures topological and aperiodicity properties of the Markov chain under consideration. This condition is used in the study of the limiting behaviour of  $(\theta_n)_{n \in \mathbb{N}}$ . Note that the condition **MD2-(ii)** is automatically satisfied adding some Gaussian noise, i.e., when  $e_\theta(X_i)$  is replaced by  $e_\theta(X_i) + p_\theta(Z_i)$  where for any  $\theta \in \Theta$ ,  $p_\theta$  is any invertible linear application from  $\mathbb{R}^d$  to  $T_\theta\Theta$  and  $(Z_i)_{i \in \mathbb{N}^*}$  is a sequence of i.i.d.  $d$ -dimensional Gaussian random variables with zero-mean and covariance matrix identity. The same assumption is considered in the Euclidean case in [Pflug \(1986, Assumption A\(v\)-A\(vi\)\)](#).

To ensure recurrence of  $(\theta_n)_{n \in \mathbb{N}}$ , we assume the existence of a Lyapunov function  $V : \Theta \rightarrow \mathbb{R}_+$  for the mean vector field  $h$ .

**H1.** (i) For any  $\theta \in \Theta$ ,  $V \circ \text{proj}_S(\theta) \leq V(\theta)$ .

(ii)  $V$  is continuously differentiable on  $\Theta$  and its Riemannian gradient  $\text{grad } V$  is geodesically  $L$ -Lipschitz, i.e., there exists  $L \geq 0$  such that for any  $\theta_0, \theta_1 \in \Theta$ , and geodesic curve  $\gamma : [0, 1] \rightarrow \Theta$  such that  $\gamma(0) = \theta_0$  and  $\gamma(1) = \theta_1$ ,

$$\|\text{grad } V(\theta_1) - T_{\theta_0}^Y \text{grad } V(\theta_0)\|_{\theta_1} \leq L\ell(\gamma), \quad (3)$$

where  $\ell(\gamma) = \|\dot{\gamma}(0)\|_{\theta_0}$  is the length of the geodesic.

(iii)  $V$  is proper on  $S$ , i.e., for any  $M \geq 0$ , there exists a compact set  $K \subset S$  such that for any  $\theta \in S \setminus K$ ,  $V(\theta) > M$ .

**H2.** There exist  $C_1 \geq 0$  and  $C_2 > 0$  such that for any  $\theta \in S$ ,  $\|h(\theta)\|_\theta^2 + C_2 \langle \text{grad } V(\theta), h(\theta) \rangle_\theta \leq C_1$ .

In addition, to quantify the convergence of  $(\theta_n)_{n \in \mathbb{N}}$  in

a neighborhood of a solution of (1), we consider the following condition for some compact set  $K^* \subset S$ .

**H3** ( $K^*$ ). *There exists  $\lambda > 0$  such that for any  $\theta \in S$ ,  $\langle \text{grad } V(\theta), h(\theta) \rangle_\theta \leq -\lambda V(\theta) \mathbb{1}_{S \setminus K^*}(\theta)$ .*

Note that under **H3**( $\emptyset$ ), if  $h(\theta) = 0$ , then  $V(\theta) = 0$  since  $V$  is a nonnegative function.

It is relevant to recognize that **H1**, **H2** and **H3** boil down to standard stability and recurrence conditions; see e.g. Benveniste et al. (1990); Duflo (1997). In the Euclidean case when we assume the uniqueness of a solution  $\theta^*$ , a common choice for  $V$  is  $\theta \mapsto \|\theta - \theta^*\|^2$  because **H1** is always satisfied. In this context, **H2** and **H3** would read  $\|h(\theta)\| \leq A + B\|\theta - \theta^*\|$  and  $\inf_{\theta \in S \setminus K^*} \langle \theta - \theta^*, h(\theta^*) \rangle / \|\theta - \theta^*\|^2 \geq \lambda$  respectively. They are therefore the Riemannian counterparts of Pflug (1986, Assumption A(i)-A(ii)). However, the square distance is no longer a suitable candidate in non-compact Riemannian settings, and therefore selecting a Lyapunov function adapted to the manifold  $\Theta$  and the geometry of the mean field  $h$  is all the more important. Note that **H1**-(iii) is automatically satisfied if  $S$  is compact. In addition, in most cases  $K^*$  and  $V$  are chosen such that  $K^* = \emptyset$  or  $\{\theta \in S : \|h(\theta)\|_\theta \leq \varepsilon\}$  for some  $\varepsilon \geq 0$ ,  $-C_2 \langle h(\theta), \text{grad } V(\theta) \rangle_\theta \geq \|h(\theta)\|_\theta^2$  for some  $C_2 > 0$  and any  $\theta \in \Theta$ , and therefore **H2** is satisfied with  $C_1 = 0$ .

The use of Lyapunov functions is really common and widespread to analyze stochastic approximation schemes, see Kushner and Yin (2003); Kushner and Huang (1981); Duflo (1997). However, compared to the Euclidean setting, the square distance cannot be used in many situations because it does not satisfy **H1**-(ii). This brought us to consider a different Lyapunov function and therefore develop an adapted framework for the Riemannian case; see Section 2.2 hereafter for more details.

We start with our first result which is established along with all the other statements of this section in the supplement Section S3.

**Theorem 1.** *Assume A1, MD1, H1-(i)-(ii), H2.*

(a) *Suppose in addition that for any  $\theta \in S$ ,  $\langle \text{grad } V(\theta), h(\theta) \rangle_\theta \leq 0$ . Then, for any  $\eta \in (0, \bar{\eta}]$ ,  $\theta_0 \in S$ , and  $n \in \mathbb{N}^*$ ,*

$$n^{-1} \sum_{k=0}^{n-1} \mathbb{E} \left[ -\langle \text{grad } V(\theta_k), h(\theta_k) \rangle_{\theta_k} \right] \leq 2V(\theta_0)/(n\eta) + \eta b, \quad (4)$$

where  $(\theta_n)_{n \in \mathbb{N}}$  is defined by (2) starting from  $\theta_0$ ,  $\bar{\eta} = [2C_2L(1 + \sigma_1^2)]^{-1}$ ,  $b = 2L\{\sigma_0^2 + C_1(1 + \sigma_1^2)\}$ .

Suppose in addition that **H3**( $K^*$ ) holds for some compact set  $K^* \subset S$ .

(b) *Then for any  $\eta \in (0, \bar{\eta}]$ ,  $\theta_0 \in S$ , and  $n \in \mathbb{N}^*$ ,*

$$n^{-1} \sum_{k=0}^{n-1} \mathbb{E}[\mathbb{1}_{S \setminus K^*}(\theta_k) V(\theta_k)] \leq V(\theta_0)/(an\eta) + \eta b/(2a), \quad (5)$$

where  $a = \lambda/2$ .

(c) *Define  $\|V\|_{K^*} = \sup\{V(\theta) : \theta \in K^*\}$  if  $K^* \neq \emptyset$  and  $\|V\|_{K^*} = 0$  otherwise. Then for any  $\eta \in (0, \bar{\eta}]$ ,  $\theta_0 \in S$ , and any  $n \in \mathbb{N}^*$ ,*

$$\mathbb{E}[V(\theta_n)] \leq \{1 - \eta a\}^n V(\theta_0) + \|V\|_{K^*} + \eta b/(2a). \quad (6)$$

Note that Theorem 1 gives, in the case  $K^* = \emptyset$ , non-asymptotic bounds of order  $\eta$  on  $n^{-1} \sum_{k=0}^{n-1} \mathbb{E}[-\langle \text{grad } V(\theta_k), h(\theta_k) \rangle_{\theta_k}]$ ,  $n^{-1} \sum_{k=0}^{n-1} \mathbb{E}[V(\theta_k)]$  and  $\mathbb{E}[V(\theta_n)]$  as  $n \rightarrow +\infty$ . In addition, the forgetting of the initial condition in (4) and (5) is linear w.r.t.  $n$ , contrary to (6) where it is exponential. A statement similar to Theorem 1-(b) holds only assuming **H1**-(i)-(ii) and replacing **H3**( $K^*$ ) by the condition that there exists  $\lambda > 0$  such that for any  $\theta \in S$ ,  $\langle \text{grad } V(\theta), h(\theta) \rangle_\theta \leq -\lambda \|h(\theta)\|_\theta^2 \mathbb{1}_{S \setminus K^*}(\theta)$ . This result is postponed to the supplement Theorem S2-Section S3.2. Theorem 1-(a) is a generalization of Hosseini and Sra (2019, Lemma 7) for SGD under a general Lyapunov condition and milder assumptions. Theorem 1-(a) is also stated with very similar assumptions in Durmus et al. (2020, Theorem 2) – except that  $\langle \text{grad } V(\theta), h(\theta) \rangle_\theta \leq 0$  in our setting is replaced with  $\|\text{grad } V(\theta)\|_\theta \leq \bar{c} \|h(\theta)\|_\theta$  – where the result is not restricted to constant step-size settings. However, Theorem 1-(c) can only be obtained for constant step-size schemes, and the following results use it as a stepping stone. We show in Section 4, how this generalization can be applied to SGD to obtain better convergence guarantees. Finally, in the same Section, we show that Theorem 1-(b)-(c) can be used to derive non-asymptotic convergence bounds for SGD applied to a geodesically strongly convex function, without any boundedness assumptions on  $\Theta$ .

The study of the asymptotic behavior of  $(\theta_n)_{n \in \mathbb{N}}$  is the second step towards understanding the quality of the approximation to the solution of (1). We now show, under suitable assumptions and for  $\eta \leq \bar{\eta}$  given in Theorem 1, first, that the chain is ergodic and admits a unique invariant distribution, and second, that this measure converges weakly to the Dirac measure at some point  $\theta^*$ , as the stepsize of the scheme goes to zero. In other words, the family of stationary distributions  $(\mu^\eta)_{\eta \in (0, \bar{\eta}]}$  concentrates around  $\theta^*$  as  $\eta \rightarrow 0$ . Possible approximations of  $\theta^*$  are therefore derived from sampling from  $\mu^\eta$  or taking its Riemannian barycenter, for a small enough  $\eta$ . If the sequence  $(\theta_n)_{n \in \mathbb{N}}$  is ergodic, then as  $n \rightarrow +\infty$  the marginal distributions of this Markov chain converge to  $\mu^\eta$  and can be used in turn as proxy to solve (1). A remaining question is



to provide an estimate of the approximation error as a function of the step-size  $\eta$ . This is tackled in Section 3.

**Theorem 2.** *Assume **A 1**, **MD 1**, **MD 2**, **H 1**, **H 2** and **H 3**( $K^*$ ) for some compact set  $K^* \subset S$ . Let  $\eta \in (0, \bar{\eta}]$  where  $\bar{\eta} = [2C_2L(1 + \sigma_1^2)]^{-1}$ . Then,  $(\theta_n)_{n \in \mathbb{N}}$  defined by (2) admits a unique stationary distribution  $\mu^\eta$  and is Harris-recurrent. In addition, there exist  $\rho \in [0, 1)$  and  $C \geq 0$  such that for any  $\theta_0 \in S$  and  $k \in \mathbb{N}$ ,  $|\mathbb{E}[g(\theta_n)] - \int_{\Theta} g(\theta) d\mu^\eta(\theta)| \leq C\rho^n(1 + V(\theta_0))$ , for any measurable function  $g : \Theta \rightarrow \mathbb{R}$  satisfying  $\sup_{\theta \in \Theta} \{|g|/V\} \leq 1$ .*

Taking  $n \rightarrow +\infty$  in Theorem 1-(c), we obtain by Theorem 2 that

$$\left| \int_{\Theta} g(\theta) d\mu^\eta(\theta) \right| \leq \|V\|_{K^*} + \eta b/(2a), \quad (7)$$

for any measurable function  $g : \Theta \rightarrow \mathbb{R}$  satisfying  $\sup_{\theta \in \Theta} \{|g|/V\} \leq 1$ . In the case  $\|V\|_{K^*} = 0$  (then  $V(\theta) = 0$  for any  $\theta \in K^*$ ), we get  $\int_{\Theta} V(\theta) d\mu^\eta(\theta) \leq \eta b/(2a)$ . Therefore, this result indicates that the family  $\{\mu^\eta : \eta \in (0, \bar{\eta}]\}$  concentrates in a  $\mathcal{O}(\eta)$ -neighborhood of  $K^*$  as  $\eta \rightarrow 0$ . In particular, if  $V$  admits a unique zero  $\theta^*$  which corresponds in many applications to a solution of (1), then we can expect that  $\{\mu^\eta : \eta \in (0, \bar{\eta}]\}$  converges in distribution to  $\delta_{\theta^*}$ , the Dirac measure at  $\theta^*$ , as  $\eta \rightarrow 0$ . The specific additional conditions to obtain such a result are the following.

**H 4.** *There exists  $\theta^* \in S$  such that for any  $r > 0$ , **H 3**( $\bar{B}(\theta^*, r)$ ) holds and that there exists  $c_r > 0$  satisfying for any  $\theta \in S \setminus \bar{B}(\theta^*, r)$ ,  $c_r \leq V(\theta)$ .*

Note that assuming **H 4** is weaker than assuming **H 3**( $\{\theta^*\}$ ) since in the first case the constant  $\lambda > 0$  in **H 3**( $\bar{B}(\theta^*, r)$ ) may depend on  $r$ . Pflug (1986) considers a similar assumption in the Euclidean case (see Pflug, 1986, Assumption A(i)).

As announced previously, we obtain the convergence in distribution of  $\{\mu^\eta : \eta \in (0, \bar{\eta}]\}$ .

**Theorem 3.** *Assume **A 1**, **MD 1**, **MD 2**, **H 1** and **H 2** and let  $\bar{\eta} = [2C_2L(1 + \sigma_1^2)]^{-1}$ .*

(a) *In addition suppose **H 3**( $K^*$ ) holds for some compact set  $K^* \subset S$  and that there exists  $c > 0$  such that for any  $\theta \in S \setminus K^*$ ,  $c \leq V(\theta)$ . Then  $\lim_{\eta \rightarrow 0} \mu^\eta\{K^*\} = 1$ , where  $\mu^\eta$  is the stationary distribution of  $Q_\eta$  for  $\eta \in (0, \bar{\eta}]$ .*

(b) *In addition suppose **H 4** holds. Then  $(\mu^\eta)_{\eta \in (0, \bar{\eta}]}$  converges weakly to  $\delta_{\theta^*}$ , as  $\eta \rightarrow 0$ .*

## 2.2 Two Examples of Lyapunov Functions

Having stated the main results of this section, we give two examples of Lyapunov functions  $V$  under the following setting for  $\Theta$ .

**A 2.**  *$\Theta$  is a Hadamard manifold. In addition, there exists  $\kappa > 0$  such that the sectional curvature of  $\Theta$  is bounded below by  $-\kappa^2$ .*

A classical choice of Lyapunov function on Euclidean spaces is  $\theta \mapsto \rho_\Theta^2(\theta, \theta^*)$ , being both strongly convex and Lipschitz-gradient. However, this function does not satisfy **H 1**-(ii) as soon as  $\Theta$  has non-zero curvature and is non-compact. In an effort to show the capital impact of curvature and in order to obtain a valid Lyapunov function satisfying the conditions **H 1** and **H 3**( $\bar{B}(\theta^*, r)$ ) for  $r > 0$ , we now introduce the necessary assumptions and consider a truncated version of  $\theta \mapsto \rho_\Theta^2(\theta, \theta^*)$ .

Let  $H = \{\text{Exp}_\theta(tH_\theta(x)) : \theta \in S, x \in X, t \in [0, \eta]\}$  be the set of all points reached from geodesics  $\gamma : [0, 1] \rightarrow \Theta$  of the form  $\gamma(0) \in S$  and  $\dot{\gamma}(0) = \eta H_{\gamma(0)}(x)$ , for any  $x \in X$ . We assume in our next result that the closure of  $H$  is compact which is implied for example in the case where  $S$  is compact and  $(\theta, x) \mapsto H_\theta(x)$  is bounded on  $S \times X$ .

**Proposition 4.** *Assume **A 2** and that the closure  $\bar{H}$  of  $H$  is compact, denote  $D_H = \text{diam}(\bar{H})$ . Consider a smooth function  $\chi_H : \Theta \rightarrow [0, 1]$  with compact support satisfying  $\chi_H(\theta) = 1$  for any  $\theta \in \bar{H}$  and for any  $\theta \in \Theta$  such that  $\inf_{\theta' \in \bar{H}} \rho_\Theta(\theta', \theta) \geq 1$ , it holds  $\chi_H(\theta) = 0$ . Consider now  $V_2 : \Theta \rightarrow \mathbb{R}_+$  defined for any  $\theta \in \Theta$  by*

$$V_2(\theta) = \chi_H(\theta) \rho_\Theta^2(\theta^*, \theta) + (1 - \chi_H(\theta)) D_H^2.$$

Then, **H 1**-(i)-(ii) holds with  $V \leftarrow V_2$  and  $L \leftarrow C_\chi(D_H + 1)(1 + \kappa \coth(\kappa D_H))$  where  $C_\chi \geq 0$  is a constant only depending on  $\chi_H$ . Suppose in addition that there exist  $r > 0, \lambda_\rho > 0$  such that for any  $\theta \in S$ ,

$$-\langle \text{Exp}_\theta^{-1}(\theta^*), h(\theta) \rangle_\theta \leq -\lambda_\rho \rho_\Theta^2(\theta^*, \theta) \mathbb{1}_{S \setminus \bar{B}(\theta^*, r)}(\theta). \quad (8)$$

Then **H 3**( $\bar{B}(\theta^*, r)$ ) holds with  $\lambda \leftarrow \lambda_\rho$ .

Note that under the setting of Proposition 4,  $V_2(\theta) \geq c$  for any  $\theta \in S \setminus \bar{B}(\theta^*, r)$  by definition, since it is continuous. Clearly, **H 1**-(iii) does not hold for  $V_2$  if  $S$  is non-compact, since  $V_2$  is constant outside of the support of  $\chi_H$ . For this reason, and to weaken the assumptions of Proposition 4, we introduce a ‘‘Huberized’’ version of the distance to  $\theta^*$ .

**Proposition 5.** *Assume **A 2**. Let  $\delta > 0$  and consider  $V_1 : \Theta \rightarrow \mathbb{R}_+$  defined for any  $\theta \in \Theta$  by*

$$V_1(\theta) = \delta^2 \left\{ (\rho_\Theta(\theta^*, \theta)/\delta)^2 + 1 \right\}^{1/2} - \delta^2. \quad (9)$$

Then, **H 1** holds with  $V \leftarrow V_1$  and  $L \leftarrow 1 + \kappa\delta$ . Suppose in addition that there exist  $r > 0, \lambda_\rho > 0$  such that for any  $\theta \in S$ , (8) holds. Then, **H 3**( $\bar{B}(\theta^*, r)$ ) holds and  $\lambda \leftarrow \lambda_\rho$ .

The proof is an extension of Durmus et al. (2020, Lemma 17) and relies heavily on this result. This Lyapunov function is still constructed upon the distance function, but as Proposition 5 shows, it is better

suit for non-positive curvature spaces. expansion of  $V$ , which is only possible because  $V$  is smooth. Note that under the setting of Proposition 5,  $V_1(\theta) \geq c$  for any  $\theta \in \mathcal{S} \setminus \bar{B}(\theta^*, r)$  by definition since it is continuous.

It is worth mentioning that if either Proposition 4 or Proposition 5 can be applied, in order to use Theorem 1 and Theorem 2 (resp. Theorem 3-(b)) the only condition to verify (relative to the Lyapunov function) is **H2** (resp. are **H2** and **H4**).

### 3 ASYMPTOTIC EXPANSION AND LAW IN THE UNCONSTRAINED CASE

The purpose of this section is to quantify the convergence derived in Theorem 3-(b). First, we establish an asymptotic expansion for the bias  $\int_{\Theta} g(\theta) d\mu^\eta(\theta) - g(\theta^*)$  w.r.t. the step size  $\eta$  for  $g$  belonging to a certain class of smooth functions from  $\Theta$  to  $\mathbb{R}$ . Our result can be applied to SGD ( $h = -\text{grad} f$ ) and implies then an asymptotic expansion of  $\int_{\Theta} \|\text{grad} f(\theta)\|_{\theta}^2 d\mu^\eta(\theta)$ . Secondly, we establish that the convergence derived in Theorem 3-(b) occurs at a rate  $\eta^{1/2}$ , through a central limit theorem for  $(\mu^\eta)_{\eta \in (0, \bar{\eta}]}$ . These two results can be understood as a bias-variance decomposition in which both terms are of order  $\eta$  and are therefore weak counterparts of Dieuleveut et al. (2017, Proposition 3, Theorem 5), Pflug (1986, Theorem 1) in a Riemannian setting. The related proofs are postponed to the supplement Section S4.

#### 3.1 Asymptotic Expansion as $\eta \rightarrow 0$

Here, we assume that **A1-(ii)** holds,  $\Theta$  is compact and the conditions of Theorem 3-(b) hold. In addition, define the covariance tensor field  $\Sigma$  on  $\Theta$ , for any  $\theta \in \Theta$  by,

$$\Sigma(\theta) = \mathbb{E}[e_{\theta}(X_1) \otimes e_{\theta}(X_1)] . \quad (10)$$

Under appropriate conditions, letting  $n \rightarrow +\infty$  in Theorem 1-(b), and Theorem S2 in the supplement, show respectively that  $\int_{\Theta} V(\theta) d\mu^\eta(\theta)$  and  $\int_{\Theta} \|h(\theta)\|_{\theta}^2 d\mu^\eta(\theta)$  are bounded by a term of order  $\eta$ . We specify this result in the case where  $h = -\text{grad} f$  for a smooth objective function  $f : \Theta \rightarrow \mathbb{R}$ . More precisely, we establish in what follows a weak asymptotic expansion for  $\int_{\Theta} \|\text{grad} f(\theta)\|_{\theta}^2 d\mu^\eta(\theta)$ , as  $\eta \rightarrow 0$  based on the following result for which we assume:

**MD3.**  $\Sigma$  is a continuous  $(2, 0)$ -tensor field on  $\Theta$ .

Denote the contraction of a covariant 2-tensor  $F$  with a contravariant 2-tensor  $G$  on  $\Theta$  by  $[F : G]$ ; see Section S6.2 (S69)-(S70) in the supplementary for more details. For two matrices  $A, B$ ,  $[A : B]$  just corresponds to  $\text{Tr}(AB^{\top})$ , where  $\top$  denotes the transpose.

**Theorem 6.** Assume **A 1-(ii)**,  $h$  is continuous,  $h(\theta^*) = 0$  and  $\Theta$  is compact. Assume also **MD1**, **MD2**, **MD3**, **H1**, **H2** and **H4**. Let  $\bar{\eta} = [2C_2L(1 + \sigma_1^2)]^{-1}$ . Then for any  $\eta \in (0, \bar{\eta}]$  and smooth function  $g : \Theta \rightarrow \mathbb{R}$ , we have

$$- \int_{\Theta} \langle \text{grad} g(\theta), h(\theta) \rangle_{\theta} d\mu^\eta(\theta) = \frac{\eta}{2} [\text{Hess} g : \Sigma](\theta^*) + \mathcal{R}_{g,\eta} ,$$

where  $\lim_{\eta \rightarrow 0} \{|\mathcal{R}_{g,\eta}|/\eta\} = 0$ .

Applying this result to SGD, i.e.  $h = -\text{grad} f$  and  $g = f$ , we obtain that  $\int_{\Theta} \|\text{grad} f(\theta)\|_{\theta}^2 d\mu^\eta(\theta) = (\eta/2) [\text{Hess} f : \Sigma](\theta^*) + \mathcal{R}_{f,\eta}$  with  $\lim_{\eta \rightarrow 0} \{|\mathcal{R}_{f,\eta}|/\eta\} = 0$ .

#### 3.2 A Central Limit Theorem on $(\mu^\eta)_{\eta \in (0, \bar{\eta}]}$

Now, we assume both **A 1-(i)** and **A 1-(ii)**, meaning  $\mathcal{S} = \Theta$  and  $\Theta$  is a Hadamard manifold. Note that under this setting  $\text{Exp}_{\theta}^{-1} : \Theta \rightarrow T_{\theta}\Theta$  is a well defined diffeomorphism for any  $\theta \in \Theta$  by (Lee, 2019, Proposition 12.9). In addition, we assume that the other conditions of Theorem 3-(b) hold. Following the approach of Pflug (1986) in Euclidean SA, to find the asymptotic rate of convergence of the family  $(\mu^\eta)_{\eta \in (0, \bar{\eta}]}$  defined in Section 2, we establish a central limit theorem in  $T_{\theta^*}\Theta$ , for the family of pushforward measures  $(\bar{\nu}^\eta)_{\eta \in (0, \bar{\eta}]}$  defined for any  $A \in \mathcal{B}(T_{\theta^*}\Theta)$  by

$$\bar{\nu}^\eta(A) = \mu^\eta \left( \text{Exp}_{\theta^*}(\eta^{1/2}A) \right) . \quad (11)$$

It is shown in Section S4 that for any  $\eta \in (0, \bar{\eta}]$ ,  $\bar{\nu}^\eta$  is the stationary distribution of the rescaled and projected Markov chain  $(\bar{U}_n)_{n \in \mathbb{N}}$  defined for any  $n \in \mathbb{N}$  by  $\bar{U}_n = \eta^{-1/2} \text{Exp}_{\theta^*}^{-1}(\theta_n)$ . Therefore, since under **A1-(i)-(ii)**, for any  $u \in T_{\theta^*}\Theta$ ,  $\rho_{\Theta}(\theta^*, \text{Exp}_{\theta^*}(u)) = \|u\|_{\theta^*}$  by (Lee, 2019, Corollary 6.12, Proposition 12.9), showing a central limit theorem for the family  $(\bar{\nu}^\eta)_{\eta \in (0, \bar{\eta}]}$  as  $\eta \rightarrow 0$  shows that asymptotically  $(\mu^\eta)_{\eta \in (0, \bar{\eta}]}$  concentrates in regions of diameter  $\mathcal{O}(\eta^{1/2})$  around  $\theta^*$  for the Riemannian distance.

We consider the following assumptions.

**MD4.** There exist  $\varepsilon_e > 0$ ,  $\tilde{\sigma}_0^2, \tilde{\sigma}_1^2 \geq 0$  such that for any  $\theta \in \Theta$ ,  $\mathbb{E}[\|e_{\theta}(X_1)\|_{\theta}^{2+\varepsilon_e}] \leq \tilde{\sigma}_0^2 + \tilde{\sigma}_1^2 V(\theta)$ .

**H5.** There exist a linear mapping  $\mathbf{A} : T_{\theta^*}\Theta \rightarrow T_{\theta^*}\Theta$  and a map  $\mathcal{H} : \Theta \rightarrow T_{\theta^*}\Theta$ , such that for any  $\theta \in \Theta$ ,

$$h(\theta) = T_{01}^{\gamma} \left( \mathbf{A} \text{Exp}_{\theta^*}^{-1}(\theta) + \mathcal{H}(\theta) \right) , \quad (12)$$

where  $\theta^*$  is defined in **H4**,  $T_{01}^{\gamma}$  denotes parallel transport along the geodesic  $\gamma : [0, 1] \rightarrow \Theta$  with  $\gamma(0) = \theta^*$  and  $\gamma(1) = \theta$ , and  $\lim_{\theta \rightarrow \theta^*} \{\|\mathcal{H}(\theta)\|_{\theta^*} / \rho_{\Theta}(\theta^*, \theta)\} = 0$ . In addition, the eigenvalues of the matrix  $\mathbf{A}$  all have strictly negative real parts. Finally, there exists  $C_3 > 0$  such that for any  $\theta \in \Theta$ ,  $\|h(\theta)\|_{\theta} \leq C_3 \rho_{\Theta}(\theta^*, \theta)$ .

**MD4** is a strengthened version of **MD 1**. We show in Theorem S16 that (12) holds in the case  $h$  is twice continuously differentiable on  $\Theta$  with  $\mathbf{A} = \nabla h(\theta^*)$ . Note

that when  $\Theta = \mathbb{R}^d$  is equipped with its Euclidean metric, (12) translates into  $h(\theta) = \mathbf{A}(\theta - \theta^*) + \mathcal{H}(\theta)$  with  $\lim_{\theta \rightarrow \theta^*} \{\|\mathcal{H}(\theta)\| / \|\theta - \theta^*\|\} = 0$ , which is exactly the assumption in Pflug (1986, Assumption B(ii)). For ease of notation, we also denote by  $\mathbf{A}$  and  $\Sigma(\theta^*)$  the matrices associated with these two linear applications in some orthonormal basis of  $T_{\theta^*}\Theta$ . **H5** guarantees the existence and uniqueness of the solution  $\mathbf{V} \in \mathbb{R}^{d \times d}$  of the Lyapunov equation  $\mathbf{A}\mathbf{V} + \mathbf{V}\mathbf{A}^\top = \Sigma(\theta^*)$ , see (Horn and Johnson, 1994, Theorem 2.2.1).

We also assume that  $V$  can be compared to a function of the distance on  $\Theta$  which leads to the strengthening of **H4**.

**H 6.** *There exists  $\theta^*$  such that **H 3**( $\{\theta^*\}$ ) holds and there exists  $\phi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  such that for any  $\theta \in \Theta$ ,  $V(\theta) \geq \phi(\rho_\Theta(\theta^*, \theta))$  and for any  $r > 0$ ,  $\inf_{[r, +\infty)} \phi > 0$ . In addition, there exists  $\bar{a} > 0$ , such that  $\lim_{r \rightarrow +\infty} \sup_{a \leq \bar{a}} a/\phi(a^{1/2}r) = 0$ .*

Note that the assumption on the growth rate of the Lyapunov function is verified when  $V = V_1$ , considered in Proposition 5. In this case, we can take  $\phi(r) = \delta^2[1 + (r/\delta)^2]^{1/2} - \delta^2$ . Pflug (1986) only considers the Euclidean case with  $V : \theta \mapsto \|\theta - \theta^*\|^2$ , in which case **H6** boils down to  $\langle \theta - \theta^*, h(\theta) \rangle \geq \|\theta - \theta^*\|^2 \lambda$ . In addition, for this particular choice of Lyapunov function, it is equivalent to **H3**( $\emptyset$ ), as the other conditions result from the existence and uniqueness of the solution  $\theta^*$  and the Euclidean properties of the squared distance function (see Pflug, 1986, Assumption A(i)-B(iii)). However, for our study, we have to generalize this assumption to any well-suited Lyapunov function, since, as we have already discussed, taking  $V : \theta \mapsto \|\theta - \theta^*\|^2$  is no longer an option in the Riemannian case.

**Theorem 7.** *Assume **A1-(i)-(ii)**, **MD1**, **MD2**, **MD3**, **MD4**, **H1**, **H2**, **H5** and **H6** hold. Suppose in addition that  $h(\theta^*) = 0$ ,  $h$  is continuous and let  $\bar{\eta} = [2C_2L(1 + \sigma_1^2)]^{-1} \wedge (4C_3)^{-1}$ . Then, the family of distributions  $(\bar{\nu}^\eta)_{\eta \in (0, \bar{\eta}]}$ , defined by (11), converges weakly to  $N(0, \mathbf{V})$  as  $\eta \rightarrow 0$  on  $T_{\theta^*}\Theta$ , where  $\mathbf{V}$  is the unique solution to the Lyapunov equation  $\mathbf{A}\mathbf{V} + \mathbf{V}\mathbf{A}^\top = \Sigma(\theta^*)$ .*

Even though  $N(0, \mathbf{V})$  is a distribution on  $\mathbb{R}^d$ , we identify  $\mathbb{R}^d$  with  $T_{\theta^*}\Theta$  using the same orthonormal basis as before. As mentioned in Section 2, Theorem 7 complements Theorem 3 because it proves that the asymptotic rate of convergence of  $(\mu^\eta)_{\eta \in (0, \bar{\eta}]}$  to  $\delta_{\theta^*}$  is  $\eta^{1/2}$ , since  $(\bar{U}_n)_{n \in \mathbb{N}}$  is rescaled by this factor with respect to the actual SA scheme  $(\theta_n)_{n \in \mathbb{N}}$ . Finally Theorem 7 can be seen as a Riemannian counterpart of Pflug (1986, Theorem 1). In the following section, we illustrate our results on SGD.

## 4 APPLICATION TO SGD

We assume throughout this section that **A1-(i)-(ii)** holds. We apply the results of Section 2 and Section 3, to the unconstrained stochastic gradient scheme, i.e.  $(\theta_n)_{n \in \mathbb{N}}$  defined by (2) with  $h = -\text{grad } f$  and  $\mathbf{S} = \Theta$ . Proofs are postponed to the supplement, Section S5.

**Geodesically Strongly Convex and Smooth Function** First, the objective function  $f : \Theta \rightarrow \mathbb{R}$  is subject to the following assumptions.

**F1.**  *$f : \Theta \rightarrow \mathbb{R}$  is twice continuously differentiable and  $\text{grad } f$  is geodesically  $L_f$ -Lipschitz, see (3).*

**F 2.**  *$f$  is continuously differentiable on  $\Theta$  and  $\lambda_f$ -strongly geodesically convex, for some  $\lambda_f > 0$ , i.e. for any  $\theta_1, \theta_2 \in \Theta$ ,  $f(\theta_2) \geq f(\theta_1) + \langle \text{Exp}_{\theta_1}^{-1}(\theta_2), \text{grad } f(\theta_1) \rangle_{\theta_1} + \lambda_f \rho_\Theta^2(\theta_1, \theta_2)$ .*

Under **F2**,  $f$  admits a unique minimizer denoted by  $\theta^*$ . In addition, we have the following inequalities.

**Lemma 8.** *Assume **A1-(i)-(ii)** and **F2**. Then for any  $\theta \in \Theta$ , we have*

$$\begin{aligned} \|\text{grad } f(\theta)\|_\theta^2 &\geq \lambda_f(f(\theta) - f(\theta^*)) \quad \text{and}, \\ f(\theta) - f(\theta^*) &\geq \lambda_f \rho_\Theta^2(\theta, \theta^*). \end{aligned} \quad (13)$$

Under **F1** and **F2**, Lemma 8 implies that  $V(\theta) = f(\theta) - f(\theta^*)$  and  $h = -\text{grad } f$  satisfy **H1** with  $L \leftarrow L_f$ , **H2** with  $C_1 \leftarrow 0, C_2 \leftarrow 1$  and **H3**( $\emptyset$ ) with  $\lambda \leftarrow \lambda_f$ . A direct application of Theorem 1-(c) leads to the following result.

**Corollary 9.** *Assume **A1-(i)-(ii)**, **MD1**, **F1**, **F2**. Consider  $(\theta_n)_{n \in \mathbb{N}}$  defined by (2) with  $h = -\text{grad } f$ . Let  $\bar{\eta} = [2L_f(1 + \sigma_1^2)]^{-1}$  and  $\eta \in (0, \bar{\eta}]$ . For any  $\theta_0 \in \Theta$ , and  $n \in \mathbb{N}$ ,*

$$\mathbb{E}[f(\theta_n) - f(\theta^*)] \leq (1 - \eta \lambda_f / 2)^n (f(\theta_0) - f(\theta^*)) + 2\eta L_f \sigma_0^2 / \lambda_f.$$

*Then, setting  $\eta = \bar{\eta} \wedge [\varepsilon \lambda_f / \{4\sigma_0^2 L_f\}]$ , for  $\varepsilon \in (0, 1)$ , and  $n = \lceil [\log(1/\varepsilon) - \log(f(\theta_0) - f(\theta^*))] / \log(1 - \eta \lambda_f / 2) \rceil$ , we get  $\mathbb{E}[f(\theta_n) - f(\theta^*)] \leq \varepsilon$ .*

Corollary 9 shows that (2) has a computational complexity of order  $\mathcal{O}(\log(1/\varepsilon)\varepsilon^{-1})$  to minimize  $f$ , without any boundedness assumptions on  $\Theta$ , contrary to Zhang and Sra (2016). In addition, Lemma 8 also implies that **H5** and **H6** hold if  $f$  is three times continuously differentiable and therefore Theorem 7 can be applied.

**Geodesically Quasi-Convex Function with Bounded Gradient** Consider the following assumption.

**F 3.**  *$f$  is twice continuously differentiable. Further, there exists  $\tilde{\lambda}_f > 0$  such that for any  $\theta \in \Theta$ ,  $-\langle \text{Exp}_\theta^{-1}(\theta^*), \text{grad } f(\theta) \rangle_\theta \geq \tilde{\lambda}_f V_1(\theta)$ , where  $V_1$  is defined by (9) with  $\delta = 1$ . In addition, there exists*

$C_f > 0$  such that for any  $\theta \in \Theta$ ,  $\|\text{grad } f(\theta)\|_\theta^2 \leq C_f(\rho_\Theta^2(\theta^*, \theta) \wedge 1)$ .

In the Euclidean case, this assumption corresponds to weak quasi-convexity as considered in [Hardt et al. \(2019, Definition 2.1\)](#),  $\langle \text{grad } f(\theta), \theta - \theta^* \rangle \geq \tau[f(\theta) - f(\theta^*)]$  where the Lyapunov function is  $\theta \mapsto f(\theta) - f(\theta^*)$ . Note that the geodesical quasi-convexity condition **F3** is a weaker version of the usual geodesical convexity in the SGD setting (see [Zhang and Sra, 2016](#)). By introducing **F3**, we can relax the condition **F1** using the following result.

**Lemma 10.** *Assume **A2** and **F2**. Suppose in addition that  $f$  is twice continuously differentiable and there exists  $M_f > 0$  such that for any  $\theta \in \Theta$ ,  $\|\text{grad } f(\theta)\|_\theta^2 \leq M_f \rho_\Theta^2(\theta^*, \theta)$ . Let  $\tilde{f} = \{f - f(\theta^*) + 1\}^{1/2}$ . Then  $\tilde{f}$  satisfies **F3** with  $C_f \leftarrow (M_f/4)[1 \wedge \lambda_f]$  and  $\tilde{\lambda}_f \leftarrow \lambda_f/(2M_f^{1/2})$ .*

Note that the condition introduced in [Lemma 10](#) is a relaxation of the condition that  $\text{grad } f$  is geodesically Lipschitz. Indeed, by [Jost \(2005, Theorem 5.6.1\)](#), for  $\theta^* \in \Theta$ ,  $\theta \mapsto \rho_\Theta^2(\theta^*, \theta)$  satisfies the conditions of [Lemma 10](#) but its gradient is not geodesically Lipschitz. A non-asymptotic bound is now given in terms of the distance-like function, defined for any  $\theta_1, \theta_2 \in \Theta$  by

$$D_\Theta^2(\theta_1, \theta_2) = \rho_\Theta^2(\theta_1, \theta_2)/(1 + \rho_\Theta^2(\theta_1, \theta_2)). \quad (14)$$

**Proposition 11.** *Assume that **A2**, **MD1**, **F3** hold. Let  $\bar{\eta} = [(8C_f/\tilde{\lambda}_f)(1 + \kappa)(1 + \sigma_1^2)]^{-1}$  and  $\eta \in (0, \bar{\eta}]$ . Consider  $(\theta_n)_{n \in \mathbb{N}}$  defined by (2) with  $h = -\text{grad } f$  and  $S = \Theta$ . Then, for any  $\theta_0 \in \Theta$  and  $n \in \mathbb{N}^*$ ,*

$$n^{-1} \sum_{k=0}^{n-1} \mathbb{E} [D_\Theta^2(\theta^*, \theta_n)] \leq 4V_1(\theta_0)/(n\eta\tilde{\lambda}_f) + 4\eta(1 + \kappa)\sigma_0^2/\tilde{\lambda}_f,$$

where  $\kappa$  is given in [A2](#), and  $V_1$  is defined by (9) with  $\delta = 1$ .

To the authors' knowledge, such a bound is novel even in a deterministic setting.

**Application to the Riemannian Barycenter Problem** To conclude our study, we consider the problem of computing the Riemannian barycenter  $\theta^*$  of a probability distribution  $\pi$  on a Hadamard manifold  $\Theta$ . First, we look at the discrete case:

$$\pi = M_\pi^{-1} \sum_{i=1}^{M_\pi} \delta_{\bar{\theta}_i}, \quad (15)$$

where  $M_\pi \in \mathbb{N}^*$  and  $\{\bar{\theta}_i\}_{i=1}^{M_\pi} \in \Theta^{M_\pi}$ . The Riemannian barycenter  $\theta^*$  or Karcher mean of  $\pi$  ([Arnaudon et al., 2012](#)) is the unique global minimum of the function  $f_\pi : \theta \mapsto \sum_{i=1}^{M_\pi} \rho_\Theta^2(\theta, \bar{\theta}_i)/(2M_\pi)$ . By [Jost \(2005, Theorem 5.6.1\)](#),  $\text{grad } f_\pi(\theta) = -M_\pi^{-1} \sum_{i=1}^{M_\pi} \text{Exp}_\theta^{-1}(\bar{\theta}_i)$  for any  $\theta \in \Theta$  and  $f_\pi$  satisfies **F2** with  $\lambda_f = 1/2$  using [Durmus et al. \(2020, Lemma 10\)](#). Therefore, by [Lemma 10](#),

[Proposition 11](#) can be applied. In addition, we get the following result, as an application of [Proposition 4](#) and [Theorem 1-\(c\)](#).

**Proposition 12.** *Assume **A2**. Let  $\theta_\pi^*$  be the Riemannian barycenter of the probability measure  $\pi$  in (15) on the Hadamard manifold  $\Theta$ , and let  $(\theta_n)_{n \in \mathbb{N}}$  be given by  $\theta_{n+1} = \text{Exp}_{\theta_n}(\eta \text{Exp}_{\theta_n}^{-1}(X_{n+1}))$ , where  $(X_n)_{n \in \mathbb{N}^*}$  is a sequence of i.i.d. random variables with distribution  $\pi$ . Then, for any  $\eta \in (0, 1/(CL_\pi^3)]$ ,  $\theta_0 \in \Theta$  and  $n \in \mathbb{N}$ ,*

$$\mathbb{E}[\rho_\Theta^2(\theta_n, \theta_\pi^*)] \leq (1 - \eta/4)^n \rho_\Theta^2(\theta_0, \theta_\pi^*) + C\eta L_\pi D^2,$$

where  $L_\pi = (1 + D)(1 + \kappa \coth(\kappa D))$ ,  $C$  is a universal constant, and  $D = \max_{i=1, \dots, M_\pi} \rho_\Theta(\theta_0, \bar{\theta}_i)$ .

Secondly, we tackle the general case where  $\pi$  is not required to be discrete or compactly supported. In this case, the mapping that we are looking to minimize is

$$f_\pi : \theta \mapsto (1/2) \int_\Theta \rho_\Theta^2(\theta, \nu) \pi(d\nu). \quad (16)$$

The function  $f_\pi$  is well-defined and finite under the following assumption.

**MD 5.** *There exists  $\theta \in \Theta$  such that  $\int_\Theta \rho_\Theta^2(\theta, \nu) \pi(d\nu) < +\infty$ .*

Note that by the triangle inequality, **MD5** is equivalent to for any  $\theta \in \Theta$  such that  $\int_\Theta \rho_\Theta^2(\theta, \nu) \pi(d\nu) < +\infty$  and therefore  $f_\pi$  is finite. Using the Lebesgue's dominated convergence theorem and [Jost \(2005, Theorem 5.6.1\)](#), we can compute its Riemannian gradient given for any  $\theta \in \Theta$  by,  $\text{grad } f_\pi(\theta) = -\int_\Theta \text{Exp}_\theta^{-1}(\nu) \pi(d\nu)$ . Then,  $f_\pi$  satisfies **F2** with  $\lambda_f = 1/2$  and admits a unique minimizer  $\theta_\pi^*$ . However,  $\text{grad } f$  does not satisfy **F1** in general. More precisely, it fails to be geodesically Lipschitz, see [Jost \(2005, Theorem 5.6.1\)](#). In the Euclidean setting, several modifications of SGD have been suggested to rescale the gradient such as RMSProp, AdaGrad and Adam ([Geoffrey, 2014; Duchi et al., 2011; Kingma and Ba, 2017](#)). Inspired by these methods, we consider the stochastic approximation scheme (2) with  $S = \Theta$  and

$$H_\theta(X_{n+1}) = (1/2) \text{Exp}_\theta^{-1} \left( X_{n+1}^{(1)} \right) \left\{ \rho_\Theta^2(\theta, X_{n+1}^{(2)})/2 + 1 \right\}^{-1/2}, \quad (17)$$

where  $X_{n+1} = (X_{n+1}^{(1)}, X_{n+1}^{(2)})$  and  $(X_k^{(1)}, X_k^{(2)})_{k \in \mathbb{N}^*}$  is an i.i.d. sequence of pairs of independent random variables with distribution  $\pi$ . The following result establishes non-asymptotic convergence bounds for the resulting recursion.

**Theorem 13.** *Assume **A2** and **MD5**. Let  $\theta_\pi^*$  be the Riemannian barycenter of the probability measure  $\pi$ . Let  $(\theta_n)_{n \in \mathbb{N}}$  be given by (2) with  $S = \Theta$  and  $H$  defined by (17). Then, for any  $n \in \mathbb{N}$ ,*

$$n^{-1} \sum_{k=0}^{n-1} \mathbb{E} [D_\Theta^2(\theta_k, \theta_\pi^*)] \leq 4V_1(\theta_0) C_\pi^{1/2} / (\eta n) + 4\eta B_\pi,$$

where  $V_1$  is defined by (9) with  $\delta \leftarrow 1$ ,  $\theta^* \leftarrow \theta_\pi^*$ ,  $C_\pi = 1 + 2f_\pi(\theta_\pi^*)$ ,  $B_\pi = (1 + \kappa)(f_\pi(\theta_\pi^*) + 1)(f_\pi(\theta_\pi^*) + 2) C_\pi^{-1/2}$  and  $D_\Theta^2$  is defined in (14).



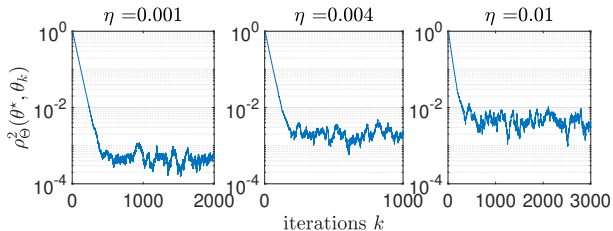


Figure 1: Paths of the algorithm in Proposition 12

## 5 NUMERICAL EXPERIMENTS

We consider in our experiments the Karcher mean estimation problem on  $\Theta = \text{Sym}_{50}^+(\mathbb{R}) \subset \mathbb{R}^{50 \times 50}$ , the symmetric definite positive matrix manifold (SPD) equipped with its affine-invariant metric, see [Pennek et al. \(2006\)](#). Note that the dimension of  $\Theta$  is 1275.

We first consider the case where  $\pi = (15)^{-1} \sum_{i=1}^{15} \delta_{x_i}$  is a discrete distribution, where  $\{x_i\}_{i=1}^{15}$  are random samples from the Wishart distribution  $\mathbf{W}(50, \text{Id})$  *i.e.* with 50 degrees of freedom and scale matrix identity. The Karcher mean  $\theta_\pi^*$  associated with  $\pi$  is estimated using the Matrix Means Toolbox ([Bini and Iannazzo, 2013](#)).

Figure 1 represents the behavior of the squared distance to the barycenter  $\theta_\pi^*$  for a single path and three step-sizes  $\eta \in \{10^{-3}, 4 \times 10^{-3}, 10^{-2}\}$ . As expected from Proposition 12, two regimes can be observed. At first, the squared-distance to the barycenter exponentially decreases and then the iterates oscillate in a  $\mathcal{O}(\eta^{1/2})$ -neighborhood of  $\theta_\pi^*$ . In addition, the rate of convergence in the exponential decay depends on the step-size.

In Figure 2, we aim at illustrating (7), Theorem 6 and Theorem 7. To this end, 1000 replications of the previous experiment are performed to obtain  $\{(\theta_n^{(i)}) : i \in \{1, \dots, 1000\}\}$  for  $n = \lceil 10/\eta \rceil$  and  $\eta \in \{1, 2.8, 4.6, 6.4, 8.2, 10\} \times 10^{-2}$ . These samples are used to estimate the mean and the variance of  $\rho_\Theta^2(\theta, \theta_\pi^*)$ , for  $\theta$  following the stationary distribution  $\mu^\eta$ . As expected, the mean and the variance are both linear w.r.t. the step-size  $\eta$ , further confirming that the iterates remain in a neighborhood of diameter  $\mathcal{O}(\eta^{1/2})$  to the ground truth.

Secondly, we examine the barycenter problem for  $\pi = \mathbf{W}(50, \text{Id})$ , following the scheme introduced in (17). The estimation of  $\theta_\pi^*$ , relative to the new distribution  $\pi$ , is now done with a 100-batch-size version of our methodology, with  $10^6$  iterations and  $\eta = 10^{-4}$ .

As a counterpart to Figure 1, in Figure 3 we are interested in the mean values of  $(D_\Theta^2(\theta_n, \theta_\pi^*))_{n \in \mathbb{N}}$  along a single path for three step-sizes  $\eta \in \{10^{-3}, 4 \times 10^{-3}, 10^{-2}\}$ , with respective burn-ins  $\{13, 3.3, 1.645\} \times 10^3$ . As pre-

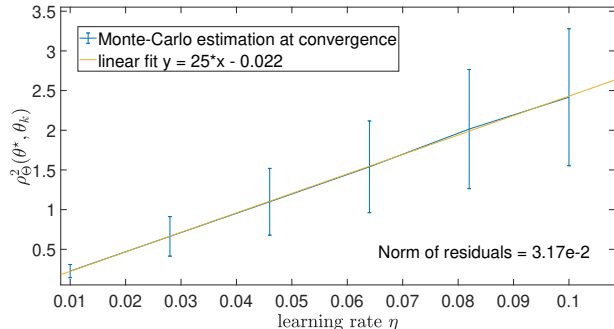


Figure 2

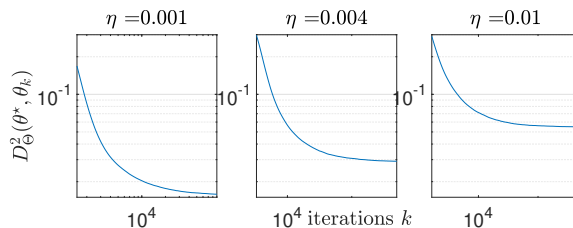


Figure 3: Paths of the algorithm in Theorem 13

dicted by Theorem 13, an initial decrease in  $\mathcal{O}(n^{-1})$  is followed by a plateau in  $\mathcal{O}(\eta)$ . We can observe that compared to Figure 1, averaging smoothes oscillations.

Finally, we also perform the experiment corresponding to Figure 2 for the discrete setting to illustrate numerically that the conclusions of (7), Theorem 6 and Theorem 7 still hold. However, due to space constraints and since the conclusions are the same than for Figure 2, the corresponding figure is postponed to the supplement Figure S1.

## Acknowledgments

AD and EM acknowledge support of the Lagrange Mathematical and Computing Research Center.

## References

- Alimisis, F., Orvieto, A., Becigneul, G., and Lucchi, A. (2020). A Continuous-time Perspective for Modeling Acceleration in Riemannian Optimization. volume 108 of *Proceedings of Machine Learning Research*, pages 1297–1307, Online. PMLR.
- Arnaudon, M., Dombry, C., Phan, A., and Yang, L. (2012). Stochastic algorithms for computing means of probability measures. *Stochastic Processes and their Applications*, 122(4):1437 – 1455.
- Bach, F. (2020). On the effectiveness of richardson extrapolation in machine learning.
- Bach, F. and Moulines, E. (2011). Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain*, pages 451–459.
- Baxter, J. and Bartlett, P. L. (2001). Infinite-horizon policy-gradient estimation. *J. Artif. Int. Res.*, 15(1):319–350.
- Benveniste, A., Métivier, M., and Priouret, P. (1990). *Adaptive algorithms and stochastic approximations*, volume 22 of *Applications of Mathematics (New York)*. Springer-Verlag, Berlin. Translated from the French by Stephen S. Wilson.
- Bini, D. A. and Iannazzo, B. (2013). Computing the Karcher mean of symmetric positive definite matrices. *Linear Algebra and its Applications*, 438(4):1700 – 1710. 16th ILAS Conference Proceedings, Pisa 2010.
- Bonnabel, S. (2013). Stochastic gradient descent on riemannian manifolds. *IEEE Transactions on Automatic Control*, 58(9):2217–2229.
- Bottou, L. (2010). Large-scale machine learning with stochastic gradient descent. In Lechevallier, Y. and Saporta, G., editors, *Proceedings of the 19th International Conference on Computational Statistics (COMPSTAT’2010)*, pages 177–187, Paris, France. Springer.
- Bottou, L. and Bousquet, O. (2008). The tradeoffs of large scale learning. In Platt, J., Koller, D., Singer, Y., and Roweis, S., editors, *Advances in Neural Information Processing Systems 20 (NIPS 2007)*, pages 161–168. NIPS Foundation (<http://books.nips.cc>).
- Boumal, N. and Absil, P.-A. (2011). RTRMC: A Riemannian trust-region method for low-rank matrix completion. In *Advances in neural information processing systems*, pages 406–414.
- Boumal, N., Voroninski, V., and Bandeira, A. (2016). The non-convex Burer-Monteiro approach works on smooth semidefinite programs. In *Advances in Neural Information Processing Systems*, pages 2757–2765.
- Cappé, O. and Moulines, E. (2009). On-line expectation-maximization algorithm for latent data models. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 71(3):593–613.
- Dieuleveut, A., Durmus, A., and Bach, F. (2017). Bridging the gap between Constant Step Size Stochastic Gradient Descent and Markov Chains.
- Duchi, J., Hazan, E., and Singer, Y. (2011). Adaptive Subgradient Methods for Online Learning and Stochastic Optimization. *Journal of Machine Learning Research*, 12(61):2121–2159.
- Duflo, M. (1997). *Random Iterative Models*. Springer-Verlag, Berlin, Heidelberg, 1st edition.
- Durmus, A., Jiménez, P., Moulines, E., Said, S., and Wai, H. (2020). Convergence analysis of Riemannian stochastic approximation schemes. *arXiv preprint arXiv:2005.13284*.
- Edelman, A., Arias, T., and Smith, S. (1998). The geometry of Algorithms with Orthogonality Constraints. *SIAM journal on Matrix Analysis and Applications*, 20(2):303–353.
- Fort, J.-C. and Pagès, G. (1999). Asymptotic behavior of a Markovian Stochastic Algorithm with Constant Step. *SIAM Journal on Control and Optimization*, 37(5):1456–1482.
- Geoffrey, H. (2014). Lecture 6e RMSprop: Divide the gradient by a running average of its recent magnitude.
- Han, A. and Gao, J. (2020). Variance reduction for Riemannian non-convex optimization with batch size adaptation.
- Hardt, M., Ma, T., and Recht, B. (2019). Gradient descent learns linear dynamical systems.
- Horn, R. A. and Johnson, C. R. (1994). *Topics in matrix analysis*. Cambridge university press.
- Hosseini, R. and Sra, S. (2019). An alternative to em for gaussian mixture models: Batch and stochastic riemannian optimization. *Mathematical Programming*, pages 1–37.
- Iannazzo, B. and Porcelli, M. (2018). The riemannian barzilai–borwein method with nonmonotone line search and the matrix geometric mean computation. *Ima Journal of Numerical Analysis*, 38:495–517.

- Ishteva, M., Absil, P.-A., Van Huffel, S., and De Lathauwer, L. (2011). Best low multilinear rank approximation of higher-order tensors, based on the riemannian trust-region scheme. *SIAM Journal on Matrix Analysis and Applications*, 32(1):115–135.
- Jaakkola, T., Jordan, M. I., and Singh, S. P. (1993). Convergence of stochastic iterative dynamic programming algorithms. In *Proceedings of the 6th International Conference on Neural Information Processing Systems*, NIPS’93, pages 703–710, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Jost, J. (2005). *Riemannian Geometry and Geometric Analysis*. Springer Universitat texts. Springer.
- Khuzani, M. B. and Li, N. (2017). Stochastic primal-dual method on riemannian manifolds with bounded sectional curvature. *arXiv preprint arXiv:1703.08167*.
- Kingma, D. P. and Ba, J. (2017). Adam: A Method for Stochastic Optimization.
- Kushner, H. J. and Huang, H. (1981). Asymptotic Properties of Stochastic Approximations with Constant Coefficients. *SIAM Journal on Control and Optimization*, 19(1):87–105.
- Kushner, H. J. and Yin, G. G. (2003). *Stochastic approximation and recursive algorithms and applications*, volume 35 of *Applications of Mathematics (New York)*. Springer-Verlag, New York, second edition. Stochastic Modelling and Applied Probability.
- Le, H. (2004). Estimation of riemannian barycentres. *Lms Journal of Computation and Mathematics*, 7:193–200.
- Lee, J. (2019). *Introduction to Riemannian Manifolds*. Springer International Publishing.
- Ma, S., Bassily, R., and Belkin, M. (2018). The power of interpolation: Understanding the effectiveness of sgd in modern over-parameterized learning. In *International Conference on Machine Learning*, pages 3325–3334.
- Nedić, A. and Bertsekas, D. (2001). *Convergence Rate of Incremental Subgradient Algorithms*, pages 223–264. Springer US, Boston, MA.
- Needell, D., Ward, R., and Srebro, N. (2014). Stochastic gradient descent, weighted sampling, and the randomized Kaczmarz algorithm. In *Advances in neural information processing systems*, pages 1017–1025.
- Pennec, X., Fillard, P., and Ayache, N. (2006). A riemannian framework for tensor computing. *International Journal of Computer Vision*, 66(1):41–66.
- Pflug, G. (1986). Stochastic Minimization with Constant Step-Size: Asymptotic Laws. *SIAM Journal on Control and Optimization*, 24(4):655–666.
- Polyak, B. T. and Juditsky, A. B. (1992). Acceleration of stochastic approximation by averaging. *SIAM J. Control Optim.*, 30(4):838–855.
- Robbins, H. and Monro, S. (1951). A stochastic approximation method. *The Annals of mathematical Statistics*, 22(3):400–407.
- Said, S. and Manton, J. (2019). The riemannian barycentre as a proxy for global optimisation. In *Geometric Science of Information, GSI*, pages 657–664.
- Sato, H., Kasai, H., and Mishra, B. (2019). Riemannian stochastic variance reduced gradient algorithm with retraction and vector transport. *SIAM Journal on Optimization*, 29(2):1444–1472.
- Sturm, K. T. (2003). Probability Measures on Metric Spaces of Nonpositive Curvature. *Contemporary Mathematics*, 338.
- Sun, J., Qu, Q., and Wright, J. (2017). Complete Dictionary Recovery Over the Sphere II: Recovery by Riemannian Trust-Region Method. *IEEE Transactions on Information Theory*, 63(2):885–914.
- Tripuraneni, N., Flammarion, N., Bach, F., and Jordan, M. I. (2018). Averaging Stochastic Gradient Descent on Riemannian Manifolds. In *Conference On Learning Theory, COLT*, pages 650–687.
- Vaswani, S., Bach, F., and Schmidt, M. (2019). Fast and faster convergence of sgd for over-parameterized models and an accelerated perceptron. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1195–1204.
- Zhang, H., Reddi, S. J., and Sra, S. (2016). Riemannian SVRG: Fast stochastic optimization on Riemannian manifolds. In *Advances in Neural Information Processing Systems*, pages 4592–4600.
- Zhang, H. and Sra, S. (2016). First-order Methods for Geodesically Convex Optimization. In *Conference on Learning Theory, COLT*, pages 1617–1638.