# Accelerated Gradient Dynamics on Riemannian Manifolds: Faster Rate and Trajectory Convergence

**Tejas Natu**[*]
Mathematics and Computer Science
Saarland University
Germany

**Camille Castera**
Department of Mathematics
University of Tübingen
Germany

**Jalal Fadili**
ENSICAEN
Normandie Université
CNRS, GREYC, France

**Peter Ochs**
Mathematics and Computer Science
Saarland University
Germany

## ABSTRACT

In order to minimize a differentiable geodesically convex function, we study a second-order dynamical system on Riemannian manifolds with an asymptotically vanishing damping term of the form $\alpha/t$. For positive values of $\alpha$, convergence rates for the objective values and convergence of trajectory is derived. We emphasize the crucial role of the curvature of the manifold for the distinction of the modes of convergence. There is a clear correspondence to the results that are known in the Euclidean case. When $\alpha$ is larger than a certain constant that depends on the curvature of the manifold, we improve the convergence rate of objective values compared to the previously known rate and prove the convergence of the trajectory of the dynamical system to an element of the set of minimizers. For $\alpha$ smaller than this curvature-dependent constant, the best known sub-optimal rates for the objective values and the trajectory are transferred to the Riemannian setting. We present computational experiments that corroborate our theoretical results.

## 1 Introduction

A perspective on constrained optimization problems that has gained substantial attention is the use of intrinsic geometry of the underlying space on which the optimization problem is posed. A typical problem is of the form

$$\min_{x \in \mathcal{M}} f(x)\,, \tag{1}$$

---

[*]Corresponding author: `natu@math.uni-sb.com`

where $\mathcal{M}$ is a Riemannian manifold and $f \colon \mathcal{M} \to \mathbb{R}$ is geodesically convex. For the special case when $\mathcal{M} = \mathbb{R}^n$, (1) becomes a smooth convex unconstrained optimization problem. Several constrained optimization problems in different fields of science and engineering can be posed as optimization problems on Riemannian manifolds. This includes the eigenvalue problem (Golub and Van Loan, 2013), the Karcher-mean problem (Bini and Iannazzo, 2013), semidefinite programming (Burer and Monteiro, 2005), Gaussian Mixture Models (Hosseini and Sra, 2015), dictionary learning (Sun et al., 2016), matrix completion (Vandereycken, 2013) and statistical shape analysis (Ring and Wirth, 2012) among others. For a more comprehensive review, see Boumal (2023). Solving (1) is challenging in general. Our interest in this paper is in the first-order methods (those using the Riemannian gradient of $f$) to solve (1). For instance, the Riemmanian gradient descent has been proposed and studied in detail, see e.g., Udriste (1994) and Zhang and Sra (2016). As a natural progression, more recently, first-order accelerated versions have also been studied on Riemannian manifolds, see e.g., Ahn and Sra (2020); Alimisis et al. (2020, 2021). Motivated from the Euclidean setting, we aim to understand the (fast) convergence of first-order algorithms on problems posed over a manifold, such as (1). To this end, we take inspiration from Nesterov's accelerated gradient algorithm (Nesterov, 1983).

For convex optimization problems in the Euclidean setting, that is, when the objective function $f$ is a convex function on $\mathbb{R}^n$, Nesterov (1983) improved the rate of convergence of vanilla gradient descent from $O\left(\frac{1}{k}\right)$ to $O(\frac{1}{k^2})$, where $k$ denotes the iteration number. In fact, this rate was proved to be optimal among all first-order methods for convex functions with Lipschitz conitnuous gradient (Nesterov, 2018). This result is a milestone in the history of convex optimization and continues to be significant with the continually growing size and scale of practical problems. However, the analysis of Nesterov's method is non-trivial. As a result, several attempts have been made to understand the underlying mathematical structure of accelerated methods and come up with different perspectives to obtain new insights. One way of understanding acceleration is to look at the continuous-time dynamics of optimization algorithms, which provides powerful analytic tools.

Su et al. (2014) proposed the following second-order in time ordinary differential equation (dynamical system) towards understanding Nesterov's accelerated gradient algorithm in the Euclidean case

$$\ddot{X}(t) + \frac{\alpha}{t}\dot{X}(t) + \nabla f(X(t)) = 0 \,, \tag{2}$$

for $t > 0$ and $\alpha > 0$, with initial conditions $X(0) = x_0$ and $\dot{X}(0) = 0$.

In this dynamical system, $\dot{X}$ and $\ddot{X}$ denote the velocity and acceleration of the trajectory $X$ respectively. A curve $X \colon [0, \infty) \to \mathbb{R}^n$ such that $X \in C^2(0, \infty) \cap C^1[0, \infty)$ is called a solution if it satisfies (2) for $t > 0$ and the stated initial conditions. Su et al. (2014) prove the existence and uniqueness of such a solution for (2).

A proper discretization of (2) gives the Nesterov's accelerated gradient method. This system is similar to the classical spring–mass–damper system and equates force given as the product of unit mass and acceleration with the gradient of a convex potential function and the damping force proportional to the velocity and an asymptotically vanishing viscous damping coefficient $\alpha/t$. In the absence of the damping term, (2) is a conservative system and therefore the presence of damping is crucial for the system to be useful to solve an optimization problem. For (2), Su et al. (2014) prove an accelerated convergence rate of $f(X(t)) - \min f = O\left(\frac{1}{t^2}\right)$ for $\alpha \geq 3$ and this has triggered an immense follow up work in the Euclidean and Hilbertian setting.

In this vast array of works that follow Su et al. (2014), we focus on the contributions that form the basis of this work. In particular, May (2017) proves the convergence rate of objective values for (2) with $\alpha > 3$ is strictly faster than $O\left(\frac{1}{t^2}\right)$ and that $f(X(t)) - \min f = o(\frac{1}{t^2})$. In infinite dimensional real Hilbert spaces, Attouch et al. (2018) show weak convergence of the trajectory to a point in the set of minimizers of $f$ (argmin$f$) provided the latter is non-empty. Attouch and Peypouquet (2016) prove little-o convergence rate for objective values and the convergence of iterates in the discrete setting. Attouch et al. (2019) analyze (2) for the case $0 < \alpha \leq 3$. They show that the rate of convergence of objective values undergoes a phase transition and that $\alpha = 3$ is the smallest constant for which the rate of convergence of objective values is $O\left(\frac{1}{t^2}\right)$. A similar analysis was carried out by Vassilis et al. (2018) and Apidopoulos et al. (2020) in the case of differential inclusion problem modeling the FISTA algorithm and the Forward–Backward algorithm.

It is natural to ask whether the above results can be extended to (1) when $\mathcal{M}$ is a Riemannian manifold. The continuous-time dynamical systems perspective for optimization on Riemannian manifolds has been studied by Munier (2007) and Alimisis et al. (2020). Munier (2007) analyzed continuous-time dynamics of steepest descent method on Riemannian manifolds and proved convergence of the trajectory to a point in the set of minimizers for geodesically convex functions. Alimisis et al. (2020) generalized (2) to Riemannian manifolds and proved $f(X(t)) - \min_{\mathcal{M}} f = O\left(\frac{1}{t^2}\right)$ when $\alpha$ is chosen appropriately in a way that takes into account the curvature of the manifold (see Section 4 for a precise meaning).

Therefore, in line with a few related approaches we study a generalization of (2) to Riemannian manifolds proposed by Alimisis et al. (2020) which we describe in detail in Section 4. Our contributions are summarized as follows;

(i) When $\alpha$ is larger than a threshold value, we prove the rate of convergence of objective values is actually $o\left(\frac{1}{t^2}\right)$, which is strictly faster than the previously known rate $O(\frac{1}{t^2})$. In addition, we prove the convergence of the trajectory to an element in the set of minimizers argmin$_{\mathcal{M}} f$.

(ii) For $\alpha$ below the threshold value, we provide convergence rates for objective values. In the same setting, we show convergence of trajectory to an element in argmin$_{\mathcal{M}} f$ under the condition that it satisfies the strong minimization property.

(iii) We perform computational experiments that confirm our theoretical guarantees.

## 2   Related Work

In the Euclidean setting, some of the earlier works studying continuous-time dynamics for first-order accelerated methods include the works of Alvarez (2000) and Attouch et al. (2000) who study (2) with a constant damping term instead of an asymptotically vanishing damping term. Cabot et al. (2009) study (2) with a general asymptotically vanishing damping term $a(t)$ and showed that when $a(t)$ is non-integrable, the solution to the dynamical system possesses optimization properties i.e., $f(X(t)) \to \min f$. Su et al. (2014) consider $a(t) = \alpha/t$ and prove accelerated convergence rate of $\mathcal{O}\left(\frac{1}{t^2}\right)$ when $\alpha \geq 3$. For $\alpha > 3$, Attouch et al. (2018) show convergence of the trajectory to an optimal solution while May (2017) shows little-o convergence rate for objective values. Attouch et al. (2019) provide convergence results for the case $0 < \alpha \leq 3$ and prove convergence of trajectory in the case where the minimizer possesses strong minimization property. Further, in the case $\alpha = 3$, for a convex objective function, they

show convergence of trajectory to the optimal solution in one dimensional problems. More recently, Attouch and Fadili (2022) have studied the Ravine method from a dynamical systems perspective and drawn similarities with the Nesterov's accelerated gradient method. A more comprehensive survey of historical aspects and research trends related to continuous-time dynamics for achieving acceleration is provided by Attouch and Fadili (2022) and Attouch and Cabot (2017).

In the Riemannian setting, standard references on optimization algorithms on Riemannian manifolds include Absil et al. (2008) and Boumal (2023). Udriste (1994) is a standard reference for convex analysis on Riemannian manifolds while Vishnoi (2018) provides a detailed pedagogical survey of geodesically convex sets and geodesically convex functions on Riemannian manifolds. Zhang and Sra (2016) develop techniques that are used to provide convergence guarantees for gradient descent method for geodesically convex functions on Riemannian manifolds. First-order accelerated algorithms for minimizing geodesically convex functions on Riemannian manifolds have been studied by Zhang and Sra (2018); Ahn and Sra (2020); Han et al. (2023); Alimisis et al. (2021). Zhang and Sra (2018) propose a computationally tractable accelerated method for geodesically strongly convex functions by proposing a new estimate sequence and show accelerated rate of convergence locally. Ahn and Sra (2020) propose the first global accelerated algorithm on Riemannian manifolds for geodesically strongly convex functions. In particular, Zhang and Sra (2018) and Ahn and Sra (2020) develop techniques to tackle metric distortion that is inherent to the analysis of algorithms on Riemannian manifolds. Han et al. (2023) generalize the work of Scieur et al. (2016) to Riemannian manifolds and propose acceleration using extrapolation. Alimisis et al. (2021) employ momentum in combination with techniques developed by Zhang and Sra (2016) to achieve acceleration and provide accelerated convergence guarantees for geodesically convex functions.

A peculiar aspect of the analysis of accelerated first-order methods on Riemannian manifolds is the set of assumptions under which the results are proved. In general, one of the standard conditions in Riemannian optimization is to assume that the exponential map is a global diffeomorphism (see Section 3 for further details). This condition ensures that the exponential map is invertible and smooth. Additionally, in order to derive results about first-order accelerated methods in both the discrete and continuous setting, we work in a bounded subset of the manifold and in particular, we must make the assumption that the trajectories in the continuous setting or iterates in the discrete setting lie in that bounded domain, see e.g., Zhang and Sra (2018), Alimisis et al. (2020) and Alimisis et al. (2021). In other words, this means that the results are valid only for trajectories or iterates that lie within that bounded subset. This is because in curved spaces, the analysis makes use of certain comparison theorems like the Rauch comparison theorem (Petersen, 2006) that dictate the size of the domain in which we must confine our analysis.

In this work, we consider the continuous-time dynamical system approach towards understanding acceleration of first-order optimization methods on Riemannian manifolds. In particular, we close the gaps in convergence guarantees between the Euclidean and Riemannian settings for the continuous-time dynamics modelling Nesterov's acceleration. We shall work with similar assumptions and for reasons discussed above. A more detailed description follows in later sections.

# 3 Preliminaries from Riemannian Geometry

We recall some basic concepts from Riemannian geometry that we shall make references to during the course of this work. This material can be found in standard references on the subject like Tu (2017), Leonor Godinho (2014) and Boumal (2023).

**Riemannian manifolds.** A smooth manifold is a Hausdorff, second-countable topological manifold such that the chart transition maps are of class $C^\infty$. To a smooth manifold, it is possible to attach at every point $x$, a real vector space called the tangent space $T_x\mathcal{M}$. The union of all tangent spaces over the manifold $\mathcal{M}$ can be imparted a smooth manifold structure and is called the tangent bundle $T\mathcal{M} = \cup_{x\in\mathcal{M}} T_x\mathcal{M}$. A smooth vector field on the manifold is a smooth map $Z\colon \mathcal{M} \to T\mathcal{M}$. A Riemannian manifold is an ordered pair $(\mathcal{M}, \langle\cdot,\cdot\rangle)$ where $\langle\cdot,\cdot\rangle$ defines an inner product $\langle\cdot,\cdot\rangle_x$ on the tangent space $T_x\mathcal{M}$ for every $x \in \mathcal{M}$. This assignment is smooth in the sense that the map $x \mapsto \langle Y_x, Z_x\rangle_x$ is a $C^\infty$ function on $\mathcal{M}$, where $Y_x$ and $Z_x$ are tangent vectors at $x$ corresponding to $C^\infty$ vector fields $Y$ and $Z$ on the manifold. The inner product on the tangent space gives the norm of a tangent vector $Y_x$ as $\|Y_x\|_x := \sqrt{\langle Y_x, Y_x\rangle_x}$. From now on, the subscript on the inner product highlighting the point on which the inner product is evaluated will be dropped when it is clear from the context.

The Riemannian inner product allows for measurement of length of a piece-wise smooth curve $\gamma\colon [a,b] \to \mathcal{M}$ using the formula $\ell(\gamma) := \int_a^b \|\gamma'(t)\|_{\gamma(t)}$, where $\gamma'(t)$ is the velocity vector field of the curve $\gamma$. This gives rise to the notion of distance between two points $x$ and $y$ on the manifold $\mathcal{M}$ given as $d(x,y) := \inf_\gamma \ell(\gamma)$, where infimum is taken over all piecewise smooth curves from $a$ to $b$ on $\mathcal{M}$ such that $\gamma(a) = x$ and $\gamma(b) = y$. The Riemannian manifold equipped with this distance becomes a metric space. We can also define the diameter of a subset $\mathcal{C}$ of the manifold $\mathcal{M}$ as $\mathrm{diam}(\mathcal{C}) = \sup_{x,y\in\mathcal{C}} d(x,y)$.

**Geodesics and parallel transport.** Geodesics generalize the notion of straight lines on curved spaces. The differentiation of vector fields along a curve is possible via the notion of a covariant derivative $\frac{D}{dt}$ associated with the unique connection operator on a Riemannian manifold called the affine or Levi–Civita connection $\nabla$. Given a curve $\gamma$ on a manifold $\mathcal{M}$ and the corresponding velocity vector field $\gamma'$. Assuming that $\gamma$ is at least $C^2-$ smooth, we call $\gamma$ a geodesic if its velocity vector is constant, i.e. $\frac{D}{dt}\left(\frac{d\gamma}{dt}\right) = 0$. Geodesics can also be defined as the solution to the variational problem of finding the curve of shortest length between two points. It should be noted that, a curve with the least distance between two points has zero acceleration, however a curve with zero acceleration need not be curve of least distance between two points. For example, on a sphere any two points can joined by both the short and long segments of the same geodesic which is the great circle.

Parallel transport refers to transporting a tangent vector along a curve such that it remains constant along the curve. On a Riemannian manifold $\mathcal{M}$ equipped with its affine connection and the covariant derivative $\frac{D}{dt}$, for any smooth curve $\gamma$ and $v \in T_{\gamma(0)}\mathcal{M}$, there exists a unique vector field $Z$ along the curve $\gamma$ such that $\frac{D}{dt}Z = 0$ and $Z(0) = v$ (Boumal, 2023). We use the notation by Zhang and Sra (2016) to denote parallel transport by $\Gamma_x^y v$ where $v \in T_x(\mathcal{M})$ is transported to $y$, that is, $T_y(\mathcal{M})$ via the geodesic $\gamma$. Parallel transport preserves inner products, i.e. $\langle u, v\rangle_x = \langle \Gamma_x^y u, \Gamma_x^y v\rangle_y$.

The notion of parallel transport allows us to define $L$- smoothness of a function $f$ defined on a Riemannian manifold $\mathcal{M}$. A function $f$ is geodesically $L$- smooth if $\left\|\mathrm{grad} f(x) - \Gamma_y^x \mathrm{grad} f(y)\right\|_x \leq L\ell(\gamma)$ for some $L > 0$ and for all $x, y \in \mathcal{M}$ and a geodesic $\gamma$.

**Exponential and logarithmic maps.**    The exponential map denoted as $\mathrm{Exp}_x\colon T_x\mathcal{M} \to \mathcal{M}$, operates on a tangent vector $v \in T_x\mathcal{M}$ and gives a point on the manifold that lies on the unique geodesic through $x$ with initial velocity $v$. The point $\mathrm{Exp}_x(v)$ lies at a distance of $\|v\|_x$ from $x$ on the geodesic. The exponential map is not injective, however on a Riemannian manifold, we can define the radius of injectivity where the exponential map is a diffeomorphism. If the radius of injectivity is non-zero, then within this neighborhood, it is possible to define the inverse of the exponential map called as logarithmic map denoted as $\mathrm{Log}_x\colon \mathcal{M} \to T_x\mathcal{M}$. $\mathrm{Log}_x(y)$ gives the tangent vector in $T_x(\mathcal{M})$ whose exponential map would give $y$ and whose length equals the distance $d(x, y)$, i.e. $d(x, y) = \|\mathrm{Log}_x y\|_x$.

**Riemannian gradient and Riemannian Hessian.**    For a smooth function $f\colon \mathcal{M} \to \mathbb{R}$, the differential $Df(x)\colon T_x\mathcal{M} \to \mathbb{R}$ is defined as $Df(x)[v] := (f \circ \gamma)'(0)$, where $v \in T_x\mathcal{M}$ and $\gamma$ is a curve on the manifold such that $\gamma(0) = x$ and $\gamma'(0) = v$. The Riemannian gradient of $f$ is the unique vector field denoted by $\mathrm{grad}f$ on $\mathcal{M}$ such that for all $(x, v) \in T\mathcal{M}$ we have $Df(x)[v] = \langle v, \mathrm{grad}f(x)\rangle_x$. The Riemannian Hessian of $f$ at $x \in \mathcal{M}$ is a linear operator $\mathrm{Hess}f\colon T_x\mathcal{M} \to T_x\mathcal{M}$ defined as $\mathrm{Hess}f(x)[v] = \nabla_v\mathrm{grad}f$.

**Sectional curvature.**    Sectional curvature generalizes the notion of Gaussian curvature of two-dimensional surfaces to higher dimensions. Starting with any two-dimensional subspace $\Pi_p$ of the tangent space at a point $x \in \mathcal{M}$, the image of $\Pi_x$ under the exponential map locally spans a two-dimensional surface $S_{\Pi_x}$ such that $T_x S_{\Pi_x} = \Pi_x$. Then the sectional curvature denoted as $K(\Pi_x)$ associated with $\Pi_x$ is the Gaussian curvature of $S_{\Pi_x}$. Since the sectional curvature is dependent of the choice of the subspace $\Pi_p$, we work with a tight global lower bound on the sectional curvature of the manifold denoted by $K_{\min}$. A similar tight global upper bound on the sectional curvature is denoted by $K_{\max}$.

**Geodesic convexity.**    The notion of convex sets and convex functions can be generalized to Riemannian manifolds by replacing straight lines with geodesics. A subset $\mathcal{C} \subset \mathcal{M}$ is called a geodesically convex set if for every $x, y \in \mathcal{C}$, there exists a geodesic $\gamma\colon [0, 1] \to \mathcal{M}$ such that $\gamma(0) = x$ and $\gamma(1) = y$ and $\gamma(t) \in \mathcal{C}$ for $t \in [0, 1]$. Since it is not necessary to have a unique geodesic between two points on a manifold (for example, two points on a sphere are joined by two segments of the same geodesic great circle of different lengths), we can define geodesically unique convex sets. A geodesically unique convex set has one geodesic segment joining any two points in the set. A function $f : \mathcal{C} \to \mathbb{R}$ is called geodesically convex function if for any geodesic $\gamma$ with $\gamma(0) = x$ and $\gamma(1) = y$ we have $f(\gamma(t)) \leq (1 - t)f(x) + tf(y)$ for all $t \in [0, 1]$. Further, a differentiable geodesically convex function satisfies

$$f(y) \geq f(x) + \langle \mathrm{grad}f(x), \mathrm{Log}_x y\rangle_x \,, \tag{3}$$

for every $x$ and $y$ in the geodesically uniquely convex set $\mathcal{C}$. We shall refer to (3) through out this work.

## 4  Problem Setting

We consider the problem in (1) and study the following dynamical system proposed by Alimisis et al. (2020) as a generalization to (2)

$$\nabla \dot{X}(t) + \frac{\alpha}{t}\dot{X}(t) + \mathrm{grad}f(X(t)) = 0\,, \; t > 0\,, \quad X(0) = x_0 \text{ and } \dot{X}(0) = 0\,, \tag{4}$$

for $\alpha > 0$ and $x_0 \in \mathcal{M}$.

In this dynamical system, $\mathrm{grad}f$ denotes the Riemannian gradient of the objective function $f$ in (1), $\dot{X}$ denotes the velocity vector field of the trajectory $X$ and $\nabla\dot{X}$ denotes the covariant derivative of the velocity vector field that generalizes the acceleration term in (2).

**Definition 1.** *A curve $X \in C^2(0,\infty) \cap C^1[0,\infty)$ on the manifold $\mathcal{M}$ that satisfies (4) is called a solution to (4).*

Alimisis et al. (2020) prove existence of solution for (4) with $\alpha > 0$ under conditions stated in Assumption 1 below. Their existence result is stated in Proposition 4.2. However, uniqueness of solution to (4) is not guaranteed.

**Assumption 1.**

    *i) The objective function $f$ is geodesically convex and geodesically L-smooth.*

    *ii) The manifold $\mathcal{M}$ is geodesically complete.*

    *iii) The exponential map is a global diffeomorphism on the manifold $\mathcal{M}$.*

**Remark 4.1.** *For a geodesically complete Riemannian manifold $\mathcal{M}$, any two points on $\mathcal{M}$ can be joined by a geodesic. The assumption of a global diffeomorphism of the exponential map ensures that the logarithmic map is well defined. Further, it is worth mentioning that an important class of manifolds called the Hadamard manifolds satisfy these conditions (Petersen, 2006).*

**Proposition 4.2.** *(Alimisis et al., 2020) Under Assumption 1, for $\alpha > 0$, System 4 has a solution $X \colon [0,\infty) \to \mathcal{M}$.*

In this work, we perform a thorough study of the asymptotic behavior of solutions to (4). For this, the curvature of the manifold plays a key role. A major difference between the Euclidean and the Riemannian setting is that the curvature of the manifold determines the choice of $\alpha$ and hence the convergence rates. To explore this, we first discuss a crucial geometric result provided by Alimisis et al. (2020).

Let $K_{\max}$ and $K_{\min}$ be the upper and lower bounds on the sectional curvature of $\mathcal{M}$ as discussed in Section 3 and fix a diameter $D$ that satisfies the following,

$$D < \frac{\pi}{\sqrt{K_{\max}}}\,, \text{ if } K_{\max} > 0\,, \quad \text{and} \quad D < \infty\,, \text{ if } K_{\max} \leq 0\,. \tag{5}$$

Consider a subset $\mathcal{C} \subset \mathcal{M}$, such that $\mathrm{diam}(C) \leq D$ where $\mathrm{diam}(C)$ denotes the diameter of $C$ (as discussed in Section 3), a curve $X \colon I \to \mathcal{C}$, where $I \subset \mathbb{R}$ and a point $z \in \mathcal{C}$. Then $d(X(t), z)$ quantifies the distance between a point $X(t)$ on the curve $X$ and the point $z$. A key step in the analysis requires a bound on the eigenvalues of $-\mathrm{Hess}\left(-\frac{1}{2}d(X(t), z)^2\right)$, where $\mathrm{Hess}$ denotes the Riemannian Hessian. This is equivalent to an expression of the form $\left\langle -\nabla_{\dot{X}}\mathrm{grad}\left(-\frac{1}{2}d(X(t), z)^2\right), \dot{X}\right\rangle$ where $\nabla_{\dot{X}}\mathrm{grad}\left(-\frac{1}{2}d(X(t), z)^2\right)$ denotes the covariant derivative of the Riemannian gradient vector field of $-\frac{1}{2}d(X(t), z)^2$. The Riemannian gradient vector field is given as

$$\mathrm{grad}\left(-\frac{1}{2}d(X(t), z)^2\right) = \mathrm{Log}_{X(t)}z\,, \tag{6}$$

7

and a proof for (6) can be found in (Pennec, 2018; Alimisis et al., 2020).

Then Alimisis et al. (2020) provide the following bounds,

$$\sigma\left(d\left(X(t),z\right)\right)\left\|\dot{X}(t)\right\|^2 \le \left\langle \nabla_{\dot{X}(t)}\mathrm{Log}_{X(t)}z, -\dot{X}(t)\right\rangle \le \xi\left(d\left(X(t),z\right)\right)\left\|\dot{X}(t)\right\|^2, \tag{7}$$

where

$$\sigma\left(p\right) := \begin{cases} 1 & , \text{ if } K_{\max} \le 0; \\ \sqrt{K_{\max}}\, p \cot\left(\sqrt{K_{\max}}\, p\right), & \text{ if } K_{\max} > 0, \end{cases} \tag{8}$$

and

$$\xi\left(p\right) := \begin{cases} \sqrt{-K_{\min}}\, p \coth\left(\sqrt{-K_{\min}}\, p\right), & \text{ if } K_{\min} < 0; \\ 1 & , \text{ if } K_{\min} \ge 0. \end{cases} \tag{9}$$
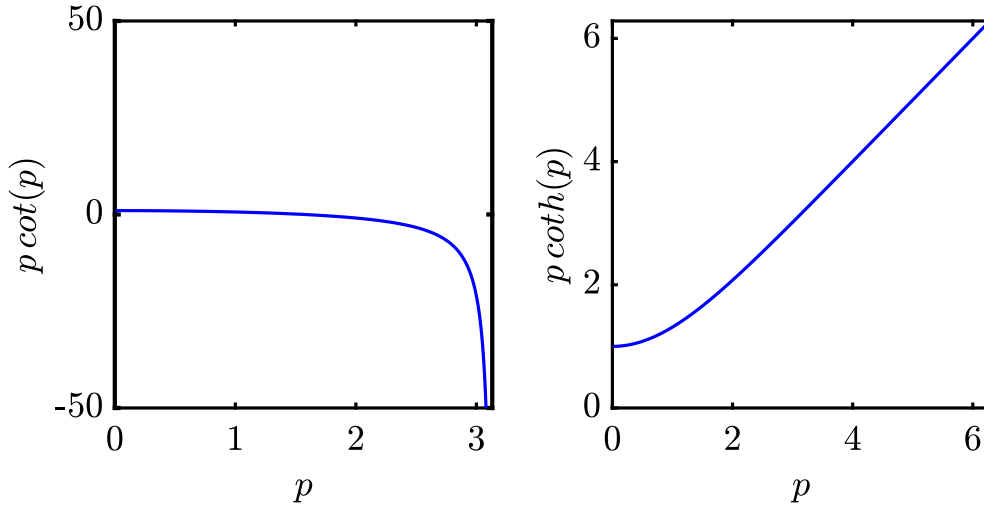


Figure 1: Functions used in the bounds given in (7) with $K_{\max} = 1$ and $K_{\min} = -1$.

The functions $p\cot(p)$ and $p\coth(p)$ are visualized in Figure 1. The terms $\sigma\left(d\left(X(t),z\right)\right)$ and $\xi\left(d\left(X(t),z\right)\right)$ are bounds on the eigenvalues of the operator $-\mathrm{Hess}\left(-\frac{1}{2}d\left(X(t),z\right)^2\right)$. From (8) and (9) we observe that the curvature of the manifold impacts these bounds. For the analysis of (4), the choice of $\alpha$ depends on the upper bound on the eigenvalues and this aspect becomes clear from the proof of Theorem 5.1 hereafter. Now, since $\xi\left(d\left(X(t),z\right)\right)$ is dependent on the parameter $t$, we instead consider an upper bound on $\xi$ by evaluating it at $D$ and define:

$$\zeta := \xi\left(D\right) \quad \text{and} \quad \delta := 2\zeta + 1. \tag{10}$$

We now make some important observations about the terms $\zeta$ and $\delta$.

(i) Since the function $p\coth(p)$ is strictly increasing for $p \in (0, \infty)$, on the set $\mathcal{C}$, we have $\xi\left(d\left(X(t),z\right)\right) \le \zeta$. Thus the upper bound in (7) can be bounded as

$$\left\langle \nabla_{\dot{X}(t)}\mathrm{Log}_{X(t)}z, -\dot{X}(t)\right\rangle \le \xi\left(d\left(X(t),z\right)\right)\left\|\dot{X}(t)\right\|^2 \le \zeta\left\|\dot{X}(t)\right\|^2. \tag{11}$$

8

(ii) When $K_{\min} \geq 0$, from (9) we have $\zeta = 1$ and thus $\delta = 3$.

(iii) When $K_{\min} < 0$, from (9), since $p \coth(\mathrm{p}) > 1$ for $p \in (0, \infty)$, we have $\zeta > 1$ and thus $\delta > 3$.

At this stage, we summarize the chain of events. We begin with a second-order system as defined in (4) and under Assumption 1, we have the existence of solutions to the system for $\alpha > 0$. We calculate $\delta$ as in (10) and this is independent of any conditions required for the existence of a solution. This is important to note because later we will analyze (4) for $0 < \alpha \leq \delta$ and $\alpha > \delta$ which is the Riemannian analog of $0 < \alpha \leq 3$ and $\alpha > 3$ in the Euclidean setting. So in the Riemannian setting, $\delta$ corresponds to the constant $3$ in the Euclidean case. Finally, to prove our main results we will make use of the bound in (11) for the case where $X$ is any solution of (4) and $z \in \operatorname{argmin}_{\mathcal{M}} f$.

Based on this discussion, we complement Assumption 1 with the following standing assumptions.

**Assumption 2.** *Let $D$ satisfy* (5).

*i) The sectional curvature of $\mathcal{M}$ is lower bounded by $K_{\min} > -\infty$.*

*ii) $\mathcal{C}$ is a geodesically convex subset of $\mathcal{M}$ with $\operatorname{diam}(\mathcal{C}) \leq D$.*

*iii) The set of minimizers $\operatorname{argmin}_{\mathcal{M}} f \neq \emptyset$ and $\operatorname{argmin}_{\mathcal{M}} f \subset \mathcal{C}$.*

*iv) The initial point $x_0 \in \mathcal{C}$ and all the solutions to* (4) *remain inside the set $\mathcal{C}$.*

**Remark 4.3.** *In order to use* (3)*, we make the assumption that $\mathcal{C}$ is geodesically convex. Since uniqueness of solution to* (4) *is not guaranteed, we make the assumption that all trajectories remain inside $\mathcal{C}$. Furthermore, we have a rather mild assumption that the set of minimizers is contained in $\mathcal{C}$. The last condition in Assumption 2 has been discussed in Section 2 as a standard assumption in the study of first-order accelerated dynamics and algorithms on Riemannian manifolds, see e.g., Zhang and Sra (2018) and Alimisis et al. (2021) (in the discrete setting) and Alimisis et al. (2020) (in the continuous-time setting).*

For (4) with $\alpha = \delta$, under Assumptions 1 and 2, Alimisis et al. (2020) prove that the convergence rate of objective values satisfies $f(X(t)) - \min_{\mathcal{M}} f = O\left(\frac{1}{t^2}\right)$. In this work, we extend the analysis by providing faster convergence rates and the convergence of solution to the set of minimizers for the case when $\alpha > \delta$. We complete the analysis by providing convergence guarantees for the case when $0 < \alpha \leq \delta$.

Thus, from the discussion in this section, we observe that the curvature of a Riemannian manifold impacts the choice of the damping coefficient $\alpha$ via a curvature-dependent term given by $\delta$. We now present the main results of this work.

## 5   Main Results

In the case $\alpha > \delta$, which is the Riemannian analog of $\alpha > 3$ in the Euclidean setting, we improve the convergence rate for objective values in Theorem 5.1 and prove the convergence of solution trajectories of (4) to an element in the set $\operatorname{argmin}_{\mathcal{M}} f$ in Theorem 5.2. In Theorem 5.3, we analyze the convergence rate in the sub-critical case $0 < \alpha \leq \delta$ and for the same setting, in Theorem 5.5, we prove convergence of trajectories under the assumption that the minimizer satisfies the strong minimization property.

9

## 5.1   Improved convergence rate when $\alpha > \delta$

Our first result improves the rate of convergence of objective values from $O\left(\frac{1}{t^2}\right)$ to $o\left(\frac{1}{t^2}\right)$. We consider (4) and perform a Lyapunov analysis similar to May (2017); Attouch et al. (2018) and prove the following result.

**Theorem 5.1.** *Assume $\alpha > \delta$ in (4). Then under Assumptions 1 and 2, any solution $X$ of (4) satisfies*

$$f(X(t)) - \min_{\mathcal{M}} f = o\left(\frac{1}{t^2}\right).$$

*Proof.* We fix some $z \in \operatorname{argmin}_{\mathcal{M}} f$, define $f^\star := \min_{\mathcal{M}} f$ and introduce the following functions $W, h \colon [t_0, \infty) \to [0, \infty)$ as

$$W(t) := \frac{1}{2}\left\|\dot{X}(t)\right\|^2 + f(X(t)) - f^\star \quad \text{and} \quad h(t) := \frac{1}{2} d\left(X(t), z\right)^2,$$

The proof strategy consists of the following steps.

(i) We show that $W'(t) \le 0$ which shows that $W(t)$ is a non-increasing function.

(ii) We show that $\lim_{t \to \infty} t^2 W(t)$ exists and is equal to some $m \ge 0$. This establishes big–O convergence rate.

(iii) We show that $\int_{t_0}^{\infty} s W(s) ds < \infty$.

(iv) Based on a simple lemma described in Appendix A.1, we must have $m = 0$.

(v) Since $W(t)$ is a sum of positive quantities, we deduce in particular that

$$\lim_{t \to \infty} t^2 \left[f\left(X(t)\right) - f^\star\right] = 0,$$

which gives us our result.

We now provide details of the proof of each step for which we will need the first and second derivatives of $h$. The derivatives of $h$ are calculated using properties of covariant derivatives of smooth vector fields on manifolds and are given as

$$h'(t) = \left\langle \operatorname{Log}_{X(t)} z, -\dot{X}(t) \right\rangle, \tag{12}$$

$$h''(t) = \left\langle \nabla_{\dot{X}(t)} \operatorname{Log}_{X(t)} z, -\dot{X}(t) \right\rangle + \left\langle \operatorname{Log}_{X(t)} z, -\nabla \dot{X}(t) \right\rangle. \tag{13}$$

A proof for (12) can be found in Alimisis et al. (2020) whereas (13) follows from the product rule for covariant derivatives, see e.g., Tu (2017)[Theorem 13.2]. An explanation of these expressions for the derivatives of $h$ including a comparison with the Euclidean case can be found in Appendix A.2.

Using the abbreviation $\kappa(t) = \frac{\alpha}{t}$, we can write

$$
\begin{aligned}
h''(t) + \kappa(t)h'(t) &= \left\langle \nabla_{\dot{X}(t)} \mathrm{Log}_{X(t)} z, -\dot{X}(t) \right\rangle + \left\langle \mathrm{Log}_{X(t)} z, -\nabla \dot{X}(t) \right\rangle + \kappa(t) \left\langle \mathrm{Log}_{X(t)} z, -\dot{X}(t) \right\rangle \\
&= \left\langle \nabla_{\dot{X}(t)} \mathrm{Log}_{X(t)} z, -\dot{X}(t) \right\rangle + \left\langle \mathrm{Log}_{X(t)} z, -\nabla \dot{X}(t) - \kappa(t)\dot{X}(t) \right\rangle \\
&= \left\langle \nabla_{\dot{X}(t)} \mathrm{Log}_{X(t)} z, -\dot{X}(t) \right\rangle + \left\langle \mathrm{Log}_{X(t)} z, \mathrm{grad} f(X(t)) \right\rangle ,
\end{aligned}
\tag{14}
$$

where the last equality follows from the definition of the ODE in (4). Next, we have

$$
\begin{aligned}
W(t) + h''(t) + \kappa(t)h'(t) = \frac{1}{2} \|X(t)\|^2 + f(X(t)) - f^\star + \left\langle \nabla_{\dot{X}(t)} \mathrm{Log}_{X(t)} z, -\dot{X}(t) \right\rangle \\
+ \left\langle \mathrm{Log}_{X(t)} z, \mathrm{grad} f(X(t)) \right\rangle .
\end{aligned}
\tag{15}
$$

From (11), we have the following curvature-dependent bound for the first term in (13),

$$
\left\langle \nabla_{\dot{X}(t)} \mathrm{Log}_{X(t)} z, -\dot{X}(t) \right\rangle \leq \zeta \left\| \dot{X}(t) \right\|^2 .
\tag{16}
$$

Using geodesic convexity of $f$ and since $f(z) = f^\star$, we can rearrange (3) as

$$
\left\langle \mathrm{Log}_{X(t)} z, \mathrm{grad} f(X(t)) \right\rangle \leq f^\star - f(X(t)) ,
$$

and combined with (16), we obtain

$$
\begin{aligned}
W(t) + h''(t) + \kappa(t)h'(t) &\leq \frac{1}{2} \left\| \dot{X}(t) \right\|^2 + (f(X(t)) - f^\star) + \zeta \left\| \dot{X}(t) \right\|^2 + (f^\star - f(X(t))) \\
&= \left( \frac{1 + 2\zeta}{2} \right) \left\| \dot{X}(t) \right\|^2 = \frac{\delta}{2} \left\| \dot{X}(t) \right\|^2 ,
\end{aligned}
\tag{17}
$$

where $\delta$ is defined in (10).

The following calculation

$$
\begin{aligned}
W'(t) &= \left\langle \nabla \dot{X}(t), \dot{X}(t) \right\rangle + \left\langle \mathrm{grad} f(X(t)), \dot{X}(t), \right\rangle \\
&= \left\langle -\kappa(t)\dot{X}(t) - \mathrm{grad} f(X(t)), \dot{X}(t) \right\rangle + \left\langle \mathrm{grad} f(X(t)), \dot{X}(t) \right\rangle \\
&= -\kappa(t) \left\langle \dot{X}(t), \dot{X}(t) \right\rangle \\
&= -\kappa(t) \left\| \dot{X}(t) \right\|^2 ,
\end{aligned}
\tag{18}
$$

shows that $W(t)$ is a non-increasing function.

Now multiply (17) by $t$, use $\kappa(t) = \frac{\alpha}{t}$ and rearrange to obtain

$$
\begin{aligned}
tW(t) &\leq t\frac{\delta}{2} \left\| \dot{X}(t) \right\|^2 - th''(t) - t\kappa(t)h'(t) \\
&= t\frac{\delta}{2} \left\| \dot{X}(t) \right\|^2 - th''(t) - \alpha h'(t) .
\end{aligned}
\tag{19}
$$

Now, $t\frac{\delta}{2}\left\|\dot{X}(t)\right\|^2$ can be written as

$$\frac{\delta}{2}t\left\|\dot{X}(t)\right\|^2 = \frac{\delta}{2}\frac{t^2}{t}\left\|\dot{X}(t)\right\|^2 = \frac{\delta}{2\alpha}t^2\kappa(t)\left\|\dot{X}(t)\right\|^2 = \frac{\delta}{2\alpha}\left(2tW(t)-(t^2W(t))'\right), \qquad (20)$$

where the last equality follows from (18). Substituting this in (19) and rearranging yields

$$\left(1-\frac{\delta}{\alpha}\right)tW(t)+\left(\frac{\delta}{2\alpha}\right)\left(t^2W(t)\right)' \le -th''(t)-\alpha h'(t). \qquad (21)$$

Integrating (21) over $[t_0,t]$, we obtain

$$\left(1-\frac{\delta}{\alpha}\right)\int_{t_0}^t sW(s)\mathrm{d}s+\left(\frac{\delta}{2\alpha}\right)(t^2W(t)) \le C_0-th'(t)+(1-\alpha)h(t), \qquad (22)$$

where $C_0 := \frac{\delta}{2\alpha}\left(t_0^2W(t_0)\right)+t_0h'(t_0)+(\alpha-1)h(t_0)$.

Using (12) and applying the Cauchy–Schwarz inequality on the tangent space at $X(t)$ provides

$$t\left|h'(t)\right| \le t\left\|\mathrm{Log}_{X(t)}z\right\|\ \left\|\dot{X}(t)\right\|. \qquad (23)$$

From the definition of $W(t)$, we have $\left\|\dot{X}(t)\right\| \le \sqrt{2W(t)}$.

Combined with observations made in (23) and using the fact that $h(t) = \frac{1}{2}\left\|\mathrm{Log}_{X(t)}z\right\|^2$ (since $d(x,y)=\left\|\mathrm{Log}_x y\right\|_x$), we obtain $t|h'(t)| \le 2\sqrt{t^2W(t)}\sqrt{h(t)}$ and therefore

$$-th'(t) \le 2\sqrt{t^2W(t)}\sqrt{h(t)}.$$

Use this in (22) to arrive at

$$\left(1-\frac{\delta}{\alpha}\right)\int_{t_0}^t sW(s)\mathrm{d}s+\left(\frac{\delta}{2\alpha}\right)(t^2W(t)) \le C_0+2\sqrt{t^2W(t)}\sqrt{h(t)}-(\alpha-1)h(t).$$

Use the inequality $-ax^2+bx \le \frac{b^2}{4a}$, $a>0$, $b\in\mathbb{R}$ with $a=(\alpha-1)$, $b=2\sqrt{t^2W(t)}$ and $x=\sqrt{h(t)}$ to obtain

$$A\int_{t_0}^t sW(s)\mathrm{d}s+Bt^2W(t) \le C_0, \qquad (24)$$

where $A := \left(1-\frac{\delta}{\alpha}\right)$ and $B := \left(\frac{\delta}{2\alpha}-\frac{1}{(\alpha-1)}\right)$.

Since $\alpha>\delta$, for both $K_{\min}\ge 0$ and $K_{\min}<0$ we have $A,B\ge 0$. Thus both the terms in (24) are non-negative and upper bounded by a constant $C_0$. Thus we infer from (24) that

$$\sup_{t\ge t_0}t^2W(t) < \infty \quad \text{and} \qquad (25)$$

$$\int_{t_0}^{+\infty} sW(s)\mathrm{d}s < \infty. \qquad (26)$$

From (20), we have

$$\left(t^2 W(t)\right)' = 2tW(t) - t^2\kappa(t)\left\|\dot{X}\right\|^2 \leq 2tW(t),$$

which combined with (26) gives us that $\int_{t_0}^{\infty}\left(t^2 W(t)\right)' < \infty$ and therefore by (25), $\lim_{t\to\infty} t^2 W(t)$ exists.

Thus, following the chain of arguments (i) − (v) stated in the beginning of the proof, we have our desired result. □

## 5.2 Convergence of trajectory when $\alpha > \delta$

We now show that for $\alpha > \delta$, the solution of (4) converges to an element in $\operatorname{argmin}_{\mathcal{M}} f$.

**Theorem 5.2.** *Assume $\alpha > \delta$ in (4). Then under Assumptions 1 and 2, there exists some $\tilde{x} \in \operatorname{argmin}_{\mathcal{M}} f$ such that $X(t) \to \tilde{x}$ as $t \to \infty$.*

*Proof.* We come back to (14). For any $z \in \operatorname{argmin}_{\mathcal{M}} f$, using geodesic convexity of $f$, we apply (3) and the fact that $f^\star - f(X(t)) \leq 0$ to (14) and obtain,

$$h''(\tau) + \kappa(\tau)h'(\tau) \leq \xi\left(d\left(X(\tau), z\right)\right)\left\|\dot{X}(\tau)\right\|^2 \leq \delta\left\|\dot{X}(\tau)\right\|^2, \tag{27}$$

since $\xi\left(d\left(X(\tau), z\right)\right) < 2\xi\left(d\left(X(\tau), z\right)\right) + 1 \leq \delta$.

Multiply both sides of (27) by $e^{\Psi(\tau,t_0)}$ where $\Psi(\tau, t_0) = \int_{t_0}^{\tau}\kappa(u)\mathrm{d}u$ is the integrating factor, and integrate from $t_0$ to $t$. Using standard integration by parts technique and the Fundamental Theorem of Calculus, we obtain,

$$h'(t) \leq e^{-\Psi(t,t_0)}h'(t_0) + \int_{t_0}^{t} e^{-\Psi(t,\tau)}\delta\left\|\dot{X}(\tau)\right\|^2 \mathrm{d}\tau, \tag{28}$$

where $-\Psi(t, \tau) = \Psi(\tau, t_0) - \Psi(t, t_0)$.

Since $\kappa(t) = \frac{\alpha}{t}$ a simple integral evaluation gives

$$\int_{s}^{\infty} e^{-\Psi(t,s)}\mathrm{d}t = \frac{s}{\alpha - 1}, \quad \forall s \geq t_0. \tag{29}$$

Further integrating (28) over $[t_0, \infty)$, we make use of (29) to obtain

$$\int_{t_0}^{\infty} h'(t)\mathrm{d}t \leq \frac{t_0}{\alpha - 1}\left|h'(t_0)\right| + \int_{t_0}^{\infty}\int_{t_0}^{t} e^{-\Psi(t,\tau)}\delta\left\|\dot{X}(\tau)\right\|^2 \mathrm{d}\tau\mathrm{d}t. \tag{30}$$

Now, upon carefully rearranging the domain of integration and subsequently applylying the Fubini's Theorem for double integrals to switch the order of integration we obtain,

$$\int_{t_0}^{\infty} h'(t)\mathrm{d}t \leq \frac{t_0}{\alpha - 1}\left|h'(t_0)\right| + \int_{t_0}^{\infty}\int_{\tau}^{\infty} e^{-\Psi(t,\tau)}\delta\left\|\dot{X}(\tau)\right\|^2 \mathrm{d}t\mathrm{d}\tau. \tag{31}$$

13

Using (29) for the inner integral in (31) we obtain,

$$\int_{t_0}^{\infty} h'(t)\mathrm{d}t \leq \frac{t_0}{\alpha-1}\,|h'(t_0)| + \frac{1}{\alpha-1}\int_{t_0}^{\infty}\tau\delta\left\|\dot{X}(\tau)\right\|^2\mathrm{d}\tau. \tag{32}$$

Finally, from (26), the right side of (32) is finite and we have that $\int_{t_0}^{\infty} h'(t)\mathrm{d}t < \infty$ which implies $\lim_{t\to\infty} h(t)$ exists. Thus we have,

$$\lim_{t\to\infty} d\left(X(t), z\right) \text{ exists for every } z \in \operatorname{argmin}_{\mathcal{M}} f. \tag{33}$$

Up to this stage, we only know that the limit in (33) exists and could be non-zero. Since the trajectory remains bounded, as a consequence of the Hopf–Rinow Theorem for complete Riemannian manifolds, there exists a subsequence $X(t_k)_{k\in\mathbb{N}}$ whose accumulation point is say $\tilde{x} \in \mathcal{C} \subset \mathcal{M}$ (Munier, 2007). By continuity of $f$, $f(\tilde{x}) = f(\lim_{k\to+\infty} X(t_k)) = \lim_{k\to+\infty} f(X(t_k)) = \lim_{t\to+\infty} f(X(t)) = f^\star$, i.e. $\tilde{x} \in \operatorname{argmin}_{\mathcal{M}} f$. Now, since (33) holds for every $z \in \operatorname{argmin}_{\mathcal{M}} f$, in particular it holds for $\tilde{x}$. This implies $d\left(X(t_k), \tilde{x}\right) \to 0$. However, by uniqueness of limit for the function, $d\left(X(t), \tilde{x}\right) \to 0$ and therefore $X(t) \to \tilde{x}$ (we emphasize that this conclusion is independent of the choice of subsequence $t_k$).

$\square$

## 5.3 Convergence rate for the sub-critical case $0 < \alpha \leq \delta$

In this section we analyze continuous-time dynamics for (4) for $0 < \alpha \leq \delta$. In the Hilbert space, this has been analyzed in Attouch et al. (2019) in the continuous-time setting while Apidopoulos et al. (2020) analyzed this in the discrete setting. We obtain a similar result for a Riemannian manifold with lower bounded sectional curvature $K_{\min}$.

**Theorem 5.3.** *Assume $0 < \alpha \leq \delta$ in (4). Then under Assumptions 1 and 2, any solution $X$ of (4) satisfies*

$$f(X(t)) - \min_{\mathcal{M}} f = O\left(\frac{1}{t^{\frac{2\alpha}{\delta}}}\right).$$

*Proof.* We fix a $z \in \operatorname{argmin}_{\mathcal{M}} f$, define $f^\star := \min_{\mathcal{M}} f$ and consider the function $W: [t_0, \infty) \to [0, \infty)$ given as

$$W(t) = A(t) + B(t) + C(t), \tag{34}$$

where

$$A(t) := t^{2p}\left(f(X(t)) - f^\star\right),$$
$$B(t) := \frac{1}{2}\left\|-\lambda(t)\left(\operatorname{Log}_{X(t)} z\right) + t^p \dot{X}(t)\right\|^2,$$
$$C(t) := = \frac{\eta(t)}{2}d\left(X(t), z\right)^2 = \frac{\eta(t)}{2}\left\|\operatorname{Log}_{X(t)} z\right\|^2,$$

and $p$ is a positive real number, $\lambda$ and $\eta$ are positive functions that will be chosen appropriately so as to make the energy function $W(t)$ non-increasing.

14

Then their derivatives are given as

$$A'(t) = 2pt^{2p-1}\left[f(X(t)) - f^\star\right] + t^{2p}\left\langle \mathrm{grad}f(X(t)), \dot{X}(t)\right\rangle,$$

$$B'(t) = \left\langle -\lambda(t)\mathrm{Log}_{X(t)}z + t^p\dot{X}(t), -\dot{\lambda}(t)\mathrm{Log}_{X(t)}z - \lambda(t)\nabla_{\dot{X}(t)}\mathrm{Log}_{X(t)}z + pt^{p-1}\dot{X}(t) + t^p\nabla\dot{X}(t)\right\rangle,$$

$$C'(t) = \frac{\dot{\eta}(t)}{2}\left\|\mathrm{Log}_{X(t)}z\right\|^2 + \eta(t)\left\langle -\dot{X}(t), \mathrm{Log}_{X(t)}z\right\rangle.$$

We make use of (4) and write $B'(t)$ as

$$B'(t) = \left\langle -\lambda(t)\mathrm{Log}_{X(t)}z + t^p\dot{X}(t), -\dot{\lambda}(t)\mathrm{Log}_{X(t)}z - \lambda(t)\nabla_{\dot{X}(t)}\mathrm{Log}_{X(t)}z + pt^{p-1}\dot{X}(t)\right.$$
$$\left. + t^{p-1}(-\alpha)\dot{X}(t) + t^p(-\mathrm{grad}f(X(t)))\right\rangle.$$

Adding the derivatives of $A$, $B$ and $C$ gives an expression for $W'(t)$. In order to avoid clutter, we avoid writing the complete expression for $W'(t)$. Instead, we make some observations about certain terms that appear in that expression that allows us to find an upper bound for $W'(t)$.

Since $d(X(t), z)^2 = \left\|\mathrm{Log}_{X(t)}z\right\|^2 = \left\langle \mathrm{Log}_{X(t)}z, \mathrm{Log}_{X(t)}z\right\rangle$, we have

$$\frac{1}{2}\frac{\mathrm{d}}{\mathrm{d}t}d\left(X(t), z\right)^2 = \frac{1}{2}\frac{\mathrm{d}}{\mathrm{d}t}\left\|\mathrm{Log}_{X(t)}z\right\| = \frac{1}{2}\frac{\mathrm{d}}{\mathrm{d}t}\left\langle \mathrm{Log}_{X(t)}z, \mathrm{Log}_{X(t)}z\right\rangle.$$

Using properties of covariant derivatives, see Tu (2017)[Theorem 13.2], we obtain,

$$\left\langle \mathrm{Log}_{X(t)}z, \nabla_{\dot{X}(t)}\mathrm{Log}_{X(t)}z\right\rangle = \frac{1}{2}\frac{\mathrm{d}}{\mathrm{d}t}d\left(X(t), z\right)^2 = \left\langle \mathrm{Log}_{X(t)}z, -\dot{X}(t)\right\rangle. \tag{35}$$

Using (3), (11) and (35) in the expression for the derivative of $W'(t)$ gives us the following bound

$$W'(t) \leq t^p\left[2pt^{p-1} - \lambda(t)\right]\left(f(X(t) - f^\star\right)$$
$$+ \left[\eta(t) + t^p\dot{\lambda}(t) - \lambda(t)(\alpha - p)t^{p-1} + \lambda(t)^2\right]\left\langle \mathrm{Log}_{X(t)}z, -\dot{X}(t)\right\rangle$$
$$- t^p\left[(\alpha - p)t^{p-1} - \lambda(t)\zeta\right]\left\|\dot{X}(t)\right\|^2 + \left[\lambda(t)\dot{\lambda}(t) + \frac{\dot{\eta}(t)}{2}\right]\left\|\mathrm{Log}_{X(t)}z\right\|^2. \tag{36}$$

We choose $\lambda$ and $\eta$ so as to make the first two terms of (36) zero. This gives

$$\lambda(t) = 2pt^{p-1} \quad \text{and} \quad \eta(t) = 2p(\alpha - 4p + 1)t^{2p-2}. \tag{37}$$

To impose that $\eta$ is non-negative, we impose the condition

$$\alpha \geq 4p - 1. \tag{38}$$

For $W$ to be a non-increasing function, we will impose the condition that $(\alpha - p)t^{p-1} - \lambda(t)\zeta \geq 0$ or equivalently

$$\alpha \geq (2\zeta + 1)p = \delta p, \tag{39}$$

15

where $\delta$ is as defined in (10). For the choice of $\lambda$ and $\eta$ as in (37), we have

$$\lambda(t)\dot{\lambda}(t) + \frac{\dot{\eta}(t)}{2} = -2p(1-p)(\alpha - 2p + 1)t^{2p-3}.$$

which is non-positive if we impose the condition

$$p \leq 1. \tag{40}$$

Thus, if we choose $p = \min\left(1, \frac{\alpha}{\delta}, \frac{\alpha+1}{4}\right)$, then conditions (38), (39) and (40) are satisfied. This implies that $W$ is a non-negative, non-increasing function associated with (4). As a consequence, we obtain

$$t^{2p}\left(f(X(t)) - f^\star\right) \leq W(t) \leq W(t_0). \tag{41}$$

Now for the case when $K_{\min} \geq 0$, by definition $\delta = 3$ and hence $0 < \alpha \leq \delta$ implies $p = \frac{\alpha}{\delta}$. For the case when $K_{\min} < 0$, we have $\delta > 3$. As a result, $0 < \alpha \leq \delta$ includes the case when $0 < \alpha < 3$ and the case when $3 \leq \alpha \leq \delta$. If $0 < \alpha < 3$ then $1 > \frac{\alpha+1}{4} > \frac{\alpha}{3} > \frac{\alpha}{\delta}$ and when $3 \leq \alpha \leq \delta$ then $\frac{\alpha+1}{4} \geq 1 \geq \frac{\alpha}{\delta}$. Thus, in general $p = \frac{\alpha}{\delta}$. Combining this with (41) we conclude the statement.

$\square$

**Corollary 5.4.** *As a corollary, we can combine this result with Theorem 5.1 and obtain the consolidated rate of convergence for $\alpha > 0$ as*

$$f(X(t)) - \min_{\mathcal{M}} f = O\left(\frac{1}{t^{p(\alpha)}}\right),$$

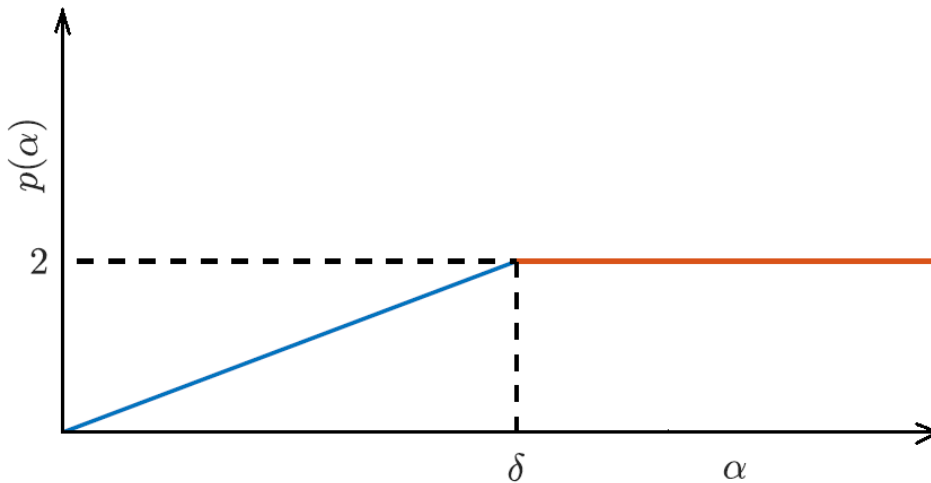*where $p(\alpha) = \min\left(2, \frac{2\alpha}{\delta}\right)$.*



Figure 2: The convergence rate in Corollary 5.4 undergoes a phase change at $\alpha = \delta$.

**Phase Transition.**   This result shows a phase transition for convergence rates at $\alpha = \delta$. For $\alpha < \delta$, the convergence rate increases linearly with a slope of $\frac{2}{\delta}$, whereas for $\alpha \geq \delta$ the convergence rate remains constant. This is shown in Figure 2 that corresponds to a similar figure in Attouch et al. (2019) in the Euclidean case that shows how $p(\alpha)$ varies as a function of $\alpha$. The rate of convergence decreases as the value of $\alpha$ decreases and this is in agreement with the previous works in the literature in the Hilbertian setting. One has then to take $\alpha$ as large as possible but the rate stagnates at $o(1/t^2)$ for $\alpha$ larger than a threshold that depends on the manifold curvature.

Additionally, we know by the work of Apidopoulos et al. (2020) that this rate is optimal for the whole space in the Hilbertian setting. It would be worth investigating whether a similar result holds for a class of manifolds with a given curvature.

### 5.4   Convergence of trajectories in the sub-critical case $0 < \alpha \leq \delta$

In the Euclidean case, when $f$ is convex, convergence of solution trajectories of (2) for the case $0 < \alpha \leq 3$ is still an open problem. However, convergence of the trajectory can be shown by assuming that the convex function has a strong minimum (Attouch et al., 2019). In the Riemannian setting we can define the notion of a strong minimum of geodesically convex function $f$ as follows.

**Definition 2.** *A geodesically convex function $f$ on a Riemannian manifold has a strong minimum if there exists $\tilde{x} \in \operatorname{argmin}_{\mathcal{M}} f$ and $\mu > 0$ such that for every $x \in \mathcal{M}$ we have*

$$f(x) \geq f(\tilde{x}) + \frac{\mu}{2} d\left(x, \tilde{x}\right)^2 . \tag{42}$$

As a result, the minimizer is actually unique. This is true in particular for geodesically strongly convex functions, for example, the Karcher-mean objective (see Section 6 for details).

**Theorem 5.5.** *Assume $0 < \alpha \leq \delta$ in (4). Then under Assumptions 1 and 2, for a geodesically convex function that admits a strong minimum $\tilde{x}$, any solution $X$ of (4) converges to $\tilde{x}$ with the rate*

$$d\left(X(t), \tilde{x}\right)^2 = O\left(\frac{1}{t^{\frac{2\alpha}{\delta}}}\right). \tag{43}$$

*Proof.* We combine Definition 2 with Theorem 5.3 to obtain

$$d\left(X(t), \tilde{x}\right)^2 \leq \frac{2}{\mu}\left(f(X(t)) - f(\tilde{x})\right),$$

and the result follows. $\square$

## 6   Numerical Experiments

We provide computational evidence for the theoretical guarantees in Section 5. In particular we would like to verify faster convergence for increasing values of $\alpha \leq \delta$ (cf. Theorem 5.3) and the little-o rate of convergence for $\alpha > \delta$ (cf. Theorem 5.1). We consider some standard optimization problems on Riemannian manifolds of positive and negative curvature. For the positive curvature, we consider the maximum eigenvalue problem and for the negative curvature we consider the Karcher-mean problem. These problems have also been considered by Sra and Hosseini (2015); Ferreira et al. (2019); Alimisis

et al. (2020). We integrate the System 4 by employing a semi-implicit discretization as in Su et al. (2014) and Alimisis et al. (2020). A description of semi-implicit discretization can be found in the Appendix B.1.

In order to demonstrate little-o rate of convergence for objective values, we study the progress of the term $t^2 (f(X(t)) - f^\star)$, where $X(t)$ is obtained from the semi-implicit solver while $f^\star$ is a benchmark value obtained from standard libraries in Matlab. For the eigenvalue problem, $f^\star$ is obtained from the Matlab eigenvalue solver whereas for the Karcher-mean problem, $f^\star$ is obtained from the Manopt library (Boumal et al., 2014). For $\alpha > \delta$, theoretically we expect $\lim_{t \to \infty} t^2 [f(X(t)) - f^\star] = 0$. Due to the limitations posed by finite machine precision, a fundamental difficulty in verifying little-o rate of convergence computationally is that the difference $f(X(t)) - f^\star$ stagnates beyond a certain stage. This allows $t^2$ to overcompensate and eventually causes their product to grow. As a result, we compute the product till the difference $f(X(t)) - f^\star$ is within a tolerance of $10^{-12}$.
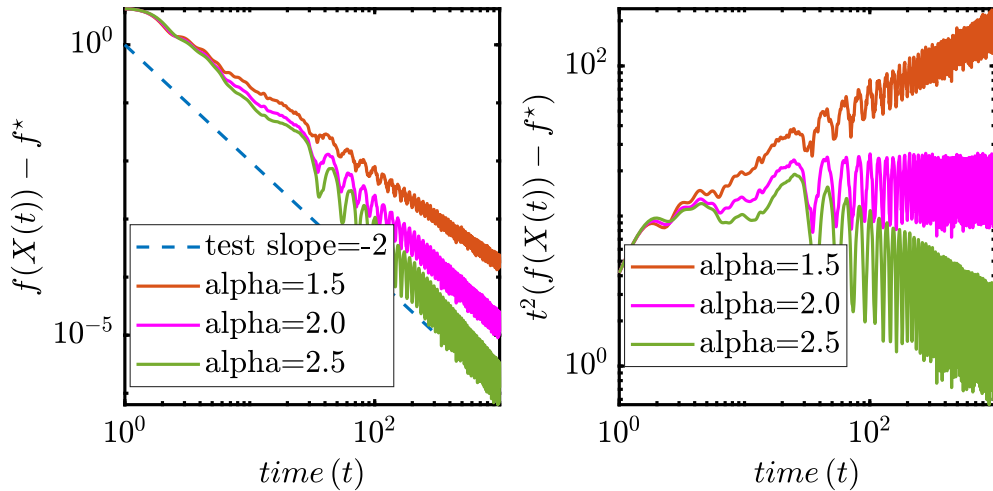


Figure 3: Convergence plots for the max–eigenvalue problem for $0 < \alpha < \delta$.

**Maximum Eigenvalue Problem.** This problem aims to find the maximum eigenvalue of a symmetric positive semi-definite matrix of large condition number. This is accomplished by minimizing the negative of the Rayleigh quotient over the hemisphere. The problem is stated as follows

$$\min_{x \in \mathbb{S}} - 0.5 x^\top A x ,$$

where $\mathbb{S} \subset \mathbb{R}^n$ is the unit hemisphere. The unit hemisphere has a constant positive curvature $K_{\min} = 1$ and hence $\delta = 3$. We refer the reader to Appendix B.2 for expressions of exponential map, Riemannian gradient and parallel transport on the sphere.

We generate the problem instance based on Alimisis et al. (2020). For the experiment, a matrix with high condition number is generated by the formula $A = \frac{1}{\beta} G^\top G$, where $G \in \mathbb{R}^{m \times n}, m > n$, is a random matrix with normally distributed entries with zero mean and variance one. We choose $m = 1000$, $n = 2500$ and $\beta = 1000$. We perform the experiment for different values of $\alpha$ with $\alpha < \delta$ and $\alpha \geq \delta$. Since $\delta = 3$, we perform experiments for $\alpha = \{1.5, 2.0, 2.5, 2.9, 3.0, 3.1, 4.0, 6.0, 8.0\}$. The system is integrated for a length of time $T = 1000$ with step size $\Delta t = 0.1$.
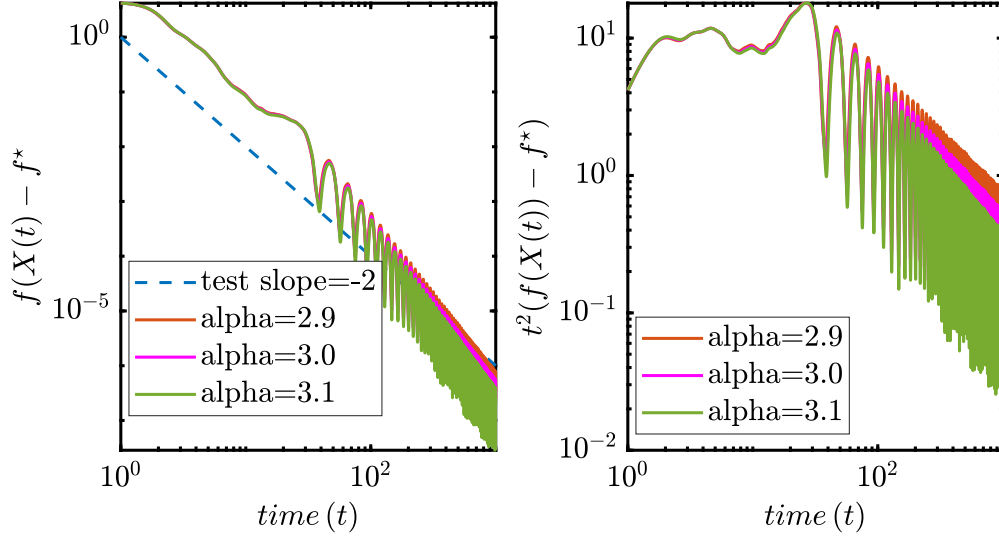
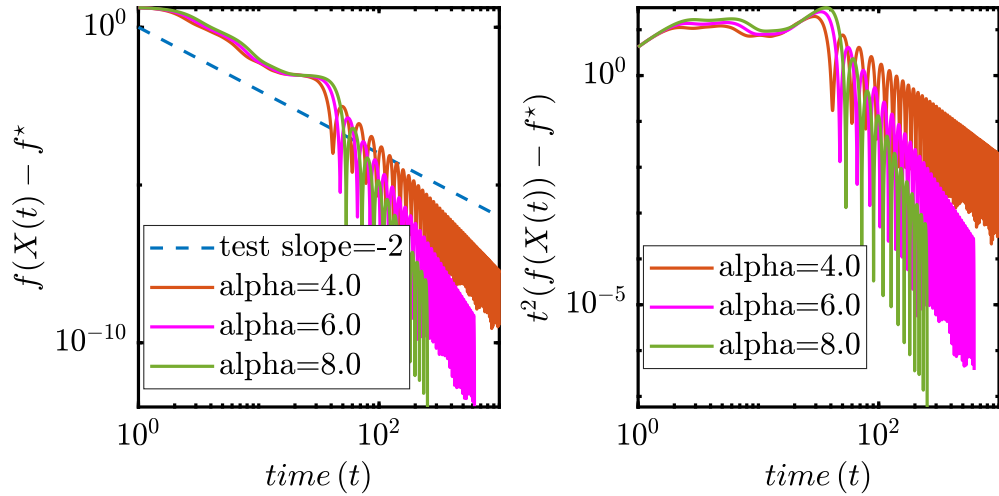Figure 4: Convergence plots for the max–eigenvalue problem–Transition.



Figure 5: Convergence plots for the max–eigenvalue problem for $\alpha > \delta$.

The results of numerical experiments are shown in Figures 3, 4 and 5 where we plot the progress of $f(X(t)) - f^\star$ and $t^2\left(f(X(t)) - f^\star\right)$ against time. While the difference $f(X(t)) - f^\star$ tends to zero for all the choices of $\alpha$, we observe an improving convergence rate for increasing values of $\alpha$. From Figure 3, it is clear that the product $t^2\left(f(X(t)) - f^\star\right)$ does not show any decay for $\alpha = 1.5$ and $\alpha = 2.0$. However, from Figure 4 we observe that close to $\delta = 3$, the product shows a decreasing trend for $\alpha = 2.5$ and $\alpha = 2.95$ and we see a transition to a convergence rate faster than $O\left(\frac{1}{t^2}\right)$. Finally, from Figure 5 for values of $\alpha > 3$ we clearly observe little-o convergence rate.

**Karcher-mean Problem.** This problem aims to find the symmetric positive definite matrix whose sum of squares of distances from a given set of symmetric positive definite matrices is the least. The problem can be posed as a Riemannian optimization problem on the manifold of symmetric positive definite matrices (SPD-manifold) with the affine–invariant metric and is a Hadamard manifold with

19

sectional curvature $K \in [-1/2, 0)$ (Criscitiello and Boumal, 2023).

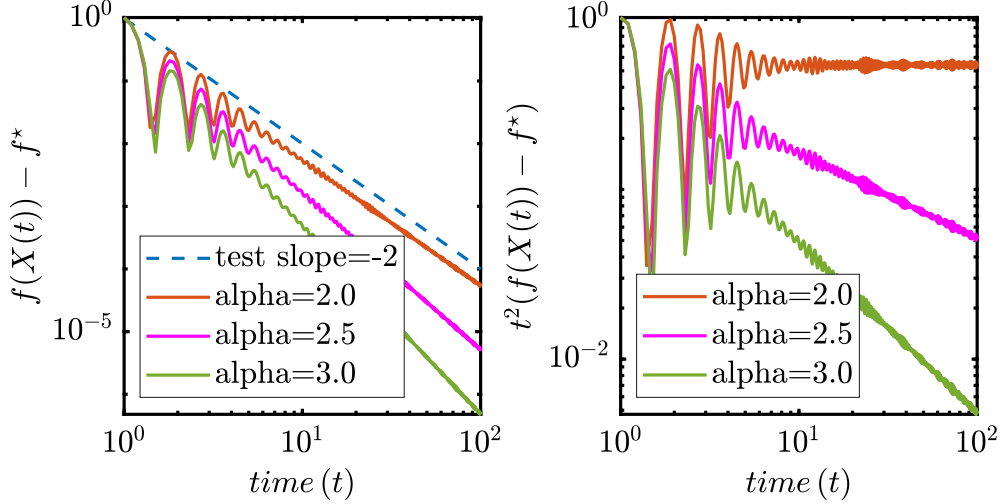$$\min_{P \in \mathbb{P}^n_{++}} \sum_{j=1}^{m} \left\| \mathrm{Logm}\left(P^{-1/2} A_j P^{-1/2}\right) \right\|_F^2$$



Figure 6: Convergence plots for the Karcher-mean problem for $0 < \alpha < \delta$.
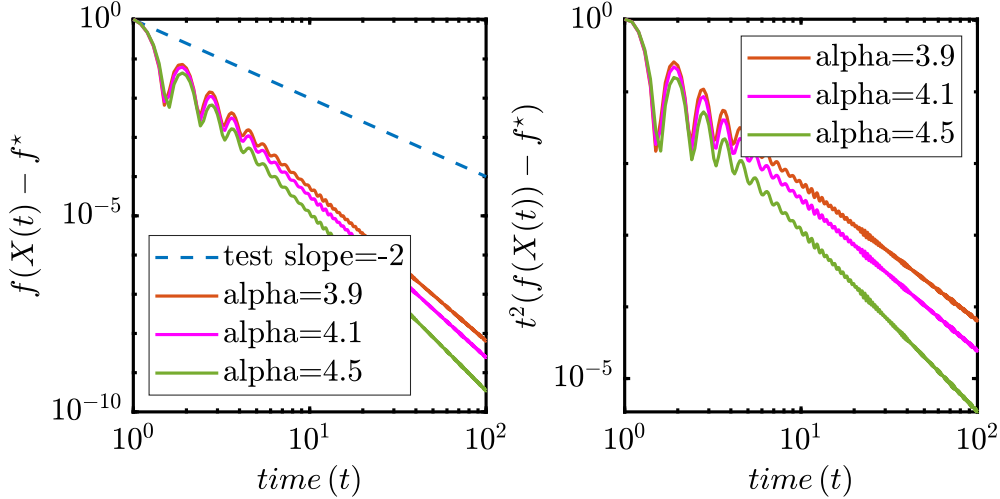


Figure 7: Convergence plots for the Karcher-mean problem–Transition.

where $\mathbb{P}^n_{++}$ denotes the SPD-manifold, $A_1, \ldots, A_m \in \mathbb{P}^n_{++}$, $\| \cdot \|_F$ denotes the Frobenius norm and $\mathrm{Logm}$ denotes matrix logarithm. The Karcher-mean problem is Euclidean non-convex problem but Riemannian strongly convex (Ferreira et al., 2019). This is an important application of Riemannian optimization where a Euclidean non-convex problem can be studied and solved as a geodesically convex problem. The expressions for the exponential map, Riemannian gradient and parallel transport have been given in Appendix B.2.
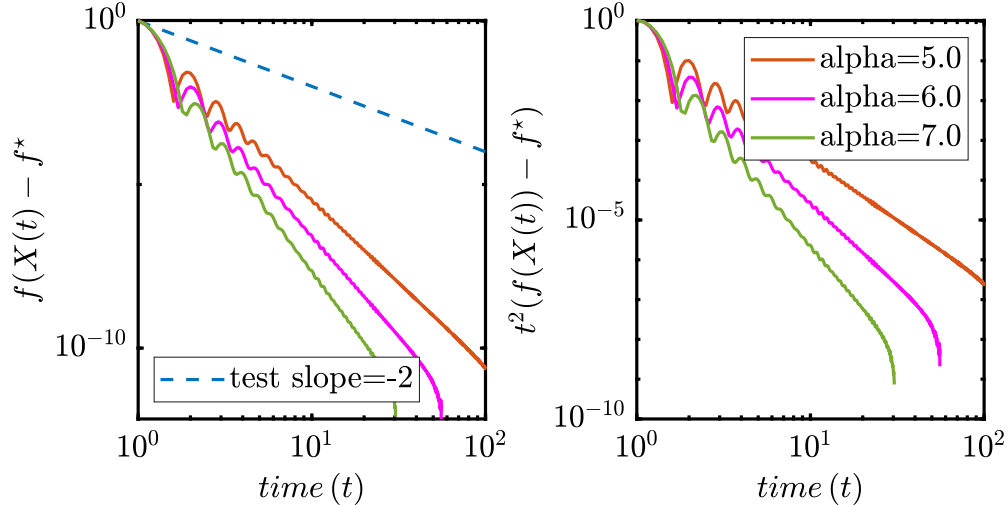
Figure 8: Convergence plots for the Karcher-mean problem for $\alpha > \delta$.

For the experiment, we compute the Karcher-mean of ten randomly generated SPD-matrices ($m = 10$) of size $n = 100$. We employ the strategy proposed in Ferreira et al. (2019) to generate SPD-matrices for the experiment and the starting point $X_0$. In order to generate the matrices, for $j = 1, \ldots, m$, we generate random orthonormal matrix $U_j$ and diagonal matrix $Q_j$ with eigenvalues in $(0, 100)$. Then $A_j = U_j Q_j U_j^\top \in \mathbb{P}_{++}^n$. The initial point $P_0$ is given as the explog–geometric mean $P_0 = \mathrm{Expm}\left(\frac{1}{m}\sum_{j=1}^m \mathrm{Logm}(A_j)\right)$, where $\mathrm{Expm}$ denotes matrix exponential.

We first estimate the value of $\delta$ for this problem. The diameter $D$ is estimated as the distance between the optimal $P \in \mathbb{P}_{++}^n$ obtained from the Manopt–solver and the initial point $X_0$. We then choose $K_{\min} = -0.1$. This gives $\zeta \approx 1.59$ and $\delta = 2\zeta + 1 \approx 4.1$. We integrate the system for a length of time $T = 100$ with step size $\Delta t = 0.1$ and for $\alpha < \delta$, $\alpha = \delta$ and $\alpha > \delta$. Since $\delta \approx 4.1$, we perform experiments for $\alpha = \{2.0, 2.5, 3.0, 3.9, 4.1, 4.5, 5.0, 6.0, 7.0\}$.

The computational experiments agree with our theoretical results. We observe that the convergence is much faster than the previous example on the sphere which is due to the fact that the Karcher-mean problem is a geodesically strongly convex. It is evident from Figures 6, 7 and 8 that while the difference $f(X(t)) - f^\star$ tends to zero for all the choices of $\alpha$, the convergence is faster for values of $\alpha$ increasing to $\delta$. While we do not observe little-o rates for $\alpha = 2.0$, we observe a decreasing trend for $t^2\left(f(X(t)) - f^\star\right)$ for $\alpha = 2.5$ and $\alpha = 3$. We observe little-o rates for values of $\alpha$ close to $\delta \approx 4.1$ and for values of $\alpha$ greater than $\delta$ as evident from Figures 7 and 8. The fact that we observe little-o rates for values of $\alpha$ close to $\delta \approx 4.1$ is due to the fact that the value of $\delta$ we have chosen is not a tight estimate and is obtained rather heuristically.

Thus, we observe that for experiments performed on manifolds of both the positive as well as the negative curvature, the computational results seem to be in line with our theoretical results.

## 7 Conclusion

In this work we studied the continuous-time dynamical system to model accelerated first-order optimization algorithms on Riemannian manifolds. We have closed gaps in the convergence guarantees

between the Euclidean setting and the Riemannian setting. In particular, corresponding to $\alpha > 3$ in the Euclidean setting, we show that the convergence rate for objective values is $o\left(\frac{1}{t^2}\right)$ for $\alpha > \delta$ on Riemannian manifolds. This rate is faster than the previously known rate of $O\left(\frac{1}{t^2}\right)$ shown by Alimisis et al. (2020). In the same setting, we also show the convergence of trajectory to an element in the set of minimizers of the objective function. We analyze the dynamical system in the sub-critical case $0 < \alpha \leq \delta$ and provide convergence rate for objective values. In this sub-critical case, we show the convergence of trajectory to a minimizer that satisfies the strong minimization property. We perform computational experiments that confirm the theoretical results.

We end this paper with some closing comments on some aspects of accelerated dynamics on Riemannian manifolds that we encountered during this work. We note that the accelerated dynamical system that we have considered cannot be studied on some rather standard manifolds like the sphere as the exponential map is not invertible on the sphere. However, this is not necessarily a limitation. This is because the sphere is a compact manifold and the only convex function that can be defined on a compact manifold is the constant function. But in general, it would be worth investigating if the assumption that the exponential map is a diffeomorphism can be relaxed.

Another important aspect is the analysis of various discretizations of the accelerated dynamical system. While we have considered the semi-implicit discretization, some other discretizations like the explicit discretization are worth considering. It would be interesting to study whether these discretizations are equivalent to the proposed first-order accelerated algorithms on Riemannian manifolds. Additionally, it is worth exploring the tightness of the value of $\delta$ as it is evident that the behavior resembling little-o rate of convergence starts to appear for values lesser than the value of $\delta$ considered. Finally, the problem of convergence of trajectories for a convex function in the case $\alpha = \delta$ is an open question even in the Euclidean case ($\alpha = 3$) for dimensions higher than one.

## Acknowledgment

# Appendices

## A    Supplementaries for Theorem 5.1

### A.1    Lemma

We draw attention to a rather simple and standard result that eventually allows us to prove the little-o rate.

**Lemma A.1.** *Consider two non-negative functions $a$ and $b$ such that $a(t)b(t)$ has a limit $m \geq 0$ as $t \to \infty$. Then, if the function b(t) is integrable and $\frac{1}{a(t)}$ is non-integrable, this means that $m = 0$.*

*Proof.* If the limit of $a(t)b(t)$ is not zero, then $b(t) \geq \frac{\tilde{m}}{a(t)}$, for some $\tilde{m} \in (0, m)$ and sufficiently large $t$, which contradicts the fact that $b(t)$ is integrable. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

### A.2    Derivatives of $h$

The expressions for the derivatives of $h(t)$ in (12) and (13) may appear rather abstract, therefore it is worthwhile to draw parallels with the Euclidean case.

Suppose $\mathcal{M} = \mathbb{R}^n$. Then the expressions for the derivatives of $h$ can be calculated by applying the chain rule and are given as

$$h'(t) = \left\langle X(t) - z, \dot{X}(t) \right\rangle, \tag{44}$$

$$h''(t) = \left\langle \dot{X}(t), \dot{X}(t) \right\rangle + \left\langle X(t) - z, \ddot{X}(t) \right\rangle = \left\| \dot{X}(t) \right\|^2 + \left\langle X(t) - z, \ddot{X}(t) \right\rangle. \tag{45}$$

In (44), the term $X(t) - z$ corresponds to the term $\mathrm{Log}_{X(t)} z$ in (12). Observe that in the Euclidean setting, $K_{\min} = K_{\max} = 0$, thus the bounds in (7) are satisfied with an equality and therefore the first term in (13) equals the first term in (45). Similarly, in the second term in (13), the covariant derivative $\nabla \dot{X}$ corresponds to the term $\ddot{X}$ in (45).

## B    Computational Results

### B.1    Semi-Implicit Discretization

In order to discretize the dynamical system, we observe that the second-order system in (4) can be written in an equivalent form as a first-order system in phase space by introducing a new variable $V$ for velocity as

$$\dot{X} = V, $$
$$\nabla V = -\frac{\alpha}{t} V - \mathrm{grad}(f(X(t))). \tag{46}$$

The semi-implicit discretization is performed by taking an explicit step in the $V$ variable using $V_k$ at the point $X_k$ to obtain $\tilde{V}_{k+1}$. The position variable is updated implicitly by applying the exponential

map on $\tilde{V}_{k+1}$ at the point $X_k$ to obtain $X_{k+1}$. Then, the updated velocity $\tilde{V}_{k+1}$ is parallel transported to $X_{k+1}$ to obtain $V_{k+1}$. For the system of differential equations (46), this is summarized as

$$\tilde{V}_{k+1} = \left(1 - \alpha\frac{\Delta t}{t_k}\right)V_k - \text{grad}(f(X_k))\Delta t\,,$$
$$X_{k+1} = \text{Exp}_{X_k}(\tilde{V}_{k+1}\Delta t)\,,$$
$$V_{k+1} = \Gamma_{X_k}^{X_{k+1}}\tilde{V}_{k+1}.$$

where $\Gamma_{X_k}^{X_{k+1}}$ denotes the parallel transport of $\tilde{V}_{k+1}$ from the point $X_k$ to $X_{k+1}$ and $\Delta t$ is the length of the time step.

## B.2 Manifold Toolbox

The expressions for exponential map, Riemannian gradient and parallel transport on a sphere can be found in Boumal (2023) or Absil et al. (2008) and are summarized as

$$\text{grad}f(x) = (I - xx^\top)(-Ax) \text{ (Riemannian Gradient)}\,,$$
$$\text{Exp}_x(v) = \cos(\|v\|)x + \sin(\|v\|)\frac{v}{\|v\|} \text{ (Exponential Map)}\,,$$
$$\Gamma_x^y(v) = v - (xx^\top)v \text{ (Parallel Transport)}\,.$$

where $I$ is an identity matrix.

The expressions for exponential map, Riemannian gradient and parallel transport on the SPD-manifold can be found in Ferreira et al. (2019); Gutman and Ho-Nguyen (2023); Axen et al. (2023) and are summarized as follows

$$\text{grad}f(P) = \sum_{i=1}^{m} X^{\frac{1}{2}}\text{Logm}(P^{\frac{1}{2}}A_i^{-1}P^{\frac{1}{2}})P^{\frac{1}{2}} \text{ (Riemannian Gradient)}$$
$$\text{Exp}_P(V) = P^{\frac{1}{2}}\text{Expm}\left(P^{-\frac{1}{2}}VP^{-\frac{1}{2}}\right)P^{\frac{1}{2}} \text{ (Exponential Map)}$$
$$\Gamma_P^{\tilde{P}}(V) = (\tilde{P}P^{-1})^{\frac{1}{2}}V(P^{-1}\tilde{P})^{\frac{1}{2}} \text{ (Parallel Transport)}.$$

where $\text{Expm}$ and $\text{Logm}$ denote the matrix exponential and logarithm.

## References

Absil, P.-A., Mahony, R., and Sepulchre, R. (2008). *Optimization algorithms on matrix manifolds.* Princeton University Press.

Ahn, K. and Sra, S. (2020). From Nesterov's estimate sequence to Riemannian acceleration. In *Conference on Learning Theory*, pages 84–118. PMLR.

Alimisis, F., Orvieto, A., Bécigneul, G., and Lucchi, A. (2020). A continuous-time perspective for modeling acceleration in Riemannian optimization. In *International Conference on Artificial Intelligence and Statistics*, pages 1297–1307. PMLR.

Alimisis, F., Orvieto, A., Becigneul, G., and Lucchi, A. (2021). Momentum improves optimization on Riemannian manifolds. In *International Conference on Artificial Intelligence and Statistics*, pages 1351–1359. PMLR.

Alvarez, F. (2000). On the minimizing property of a second order dissipative system in Hilbert spaces. *SIAM Journal on Control and Optimization*, 38(4):1102–1119.

Apidopoulos, V., Aujol, J.-F., and Dossal, C. (2020). Convergence rate of inertial forward–backward algorithm beyond Nesterov's rule. *Mathematical Programming*, 180(1-2):137–156.

Attouch, H. and Cabot, A. (2017). Asymptotic stabilization of inertial gradient dynamics with time-dependent viscosity. *Journal of Differential Equations*, 263(9):5412–5458.

Attouch, H., Chbani, Z., Peypouquet, J., and Redont, P. (2018). Fast convergence of inertial dynamics and algorithms with asymptotic vanishing viscosity. *Mathematical Programming*, 168:123–175.

Attouch, H., Chbani, Z., and Riahi, H. (2019). Rate of convergence of the Nesterov accelerated gradient method in the subcritical case $\alpha \leq 3$. *ESAIM: Control, Optimisation and Calculus of Variations*, 25:2.

Attouch, H. and Fadili, J. (2022). From the Ravine method to the Nesterov method and vice versa: a dynamical system perspective. *SIAM Journal on Optimization*, 32(3):2074–2101.

Attouch, H., Goudou, X., and Redont, P. (2000). The heavy ball with friction method, i. the continuous dynamical system: global exploration of the local minima of a real-valued function by asymptotic analysis of a dissipative dynamical system. *Communications in Contemporary Mathematics*, 2(1):1–34.

Attouch, H. and Peypouquet, J. (2016). The rate of convergence of Nesterov's accelerated forward-backward method is actually faster than $1/k^2$. *SIAM Journal on Optimization*, 26(3):1824–1834.

Axen, S. D., Baran, M., Bergmann, R., and Rzecki, K. (2023). Manifolds.jl: An extensible Julia framework for data analysis on manifolds. *AMS Transactions on Mathematical Software*. accepted for publication.

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.

Boumal, N. (2023). *An Introduction to Optimization on Smooth Manifolds*. Cambridge University Press.

Boumal, N., Mishra, B., Absil, P.-A., and Sepulchre, R. (2014). Manopt, a Matlab toolbox for optimization on manifolds. *Journal of Machine Learning Research*, 15(1):1455–1459.

Burer, S. and Monteiro, R. D. (2005). Local minima and convergence in low-rank semidefinite programming. *Mathematical programming*, 103(3):427–444.

Cabot, A., Engler, H., and Gadat, S. (2009). On the long time behavior of second order differential equations with asymptotically small dissipation. *Transactions of the American Mathematical Society*, 361(11):5983–6017.

Criscitiello, C. and Boumal, N. (2023). An accelerated first-order method for non-convex optimization on manifolds. *Foundations of Computational Mathematics*, 23(4):1433–1509.

Ferreira, O. P., Louzeiro, M. S., and Prudente, L. (2019). Gradient method for optimization on Riemannian manifolds with lower bounded curvature. *SIAM Journal on Optimization*, 29(4):2517–2541.

Golub, G. H. and Van Loan, C. F. (2013). *Matrix computations*. JHU press.

Gutman, D. H. and Ho-Nguyen, N. (2023). Coordinate descent without coordinates: Tangent subspace descent on Riemannian manifolds. *Mathematics of Operations Research*, 48(1):127–159.

Han, A., Mishra, B., Jawanpuria, P., and Gao, J. (2023). Riemannian accelerated gradient methods via extrapolation. In *International Conference on Artificial Intelligence and Statistics*, pages 1554–1585. PMLR.

Hosseini, R. and Sra, S. (2015). Matrix manifold optimization for Gaussian mixtures. *Advances In Neural Information Processing Systems*, 28.

Leonor Godinho, J. N. (2014). *An Introduction to Riemannian Geometry With Applications to Mechanics and Relativity*. Springer Cham.

May, R. (2017). Asymptotic for a second-order evolution equation with convex potential and vanishing damping term. *Turkish Journal of Mathematics*, 41(3):681–685.

Munier, J. (2007). Steepest descent method on a Riemannian manifold: the convex case. *Balkan Journal of Geometry & Its Applications*, 12(2).

Nesterov, Y. (2018). *Lectures on Convex Optimization*. Springer Cham.

Nesterov, Y. E. (1983). A method of solving a convex programming problem with convergence rate $O(1/k^2)$. In *Doklady Akademii Nauk*, volume 269, pages 543–547. Russian Academy of Sciences.

Pennec, X. (2018). Barycentric subspace analysis on manifolds.

Petersen, P. (2006). *Riemannian geometry*, volume 171. Springer.

Ring, W. and Wirth, B. (2012). Optimization methods on Riemannian manifolds and their application to shape space. *SIAM Journal on Optimization*, 22(2):596–627.

Scieur, D., d'Aspremont, A., and Bach, F. (2016). Regularized nonlinear acceleration. *Advances In Neural Information Processing Systems*, 29.

Sra, S. and Hosseini, R. (2015). Conic geometric optimization on the manifold of positive definite matrices. *SIAM Journal on Optimization*, 25(1):713–739.

Su, W., Boyd, S., and Candes, E. (2014). A differential equation for modeling Nesterov's accelerated gradient method: theory and insights. *Advances In Neural Information Processing Systems*, 27.

Sun, J., Qu, Q., and Wright, J. (2016). Complete dictionary recovery over the sphere i: Overview and the geometric picture. *IEEE Transactions on Information Theory*, 63(2):853–884.

Tu, L. W. (2017). *Differential Geometry: Connections, Curvature, and Characteristic Classes*. Springer Cham.

Udriste, C. (1994). *Convex functions and optimization methods on Riemannian manifolds*, volume 297. Springer Science & Business Media.

Vandereycken, B. (2013). Low-rank matrix completion by Riemannian optimization. *SIAM Journal on Optimization*, 23(2):1214–1236.

Vassilis, A., Jean-François, A., and Charles, D. (2018). The differential inclusion modeling FISTA algorithm and optimality of convergence rate in the case $b \leq 3$. *SIAM Journal on Optimization*, 28(1):551–574.

Vishnoi, N. K. (2018). Geodesic convex optimization: Differentiation on manifolds. *arXiv: 1806.06373*.

Zhang, H. and Sra, S. (2016). First-order methods for geodesically convex optimization. In *Conference on Learning Theory*, pages 1617–1638. PMLR.

Zhang, H. and Sra, S. (2018). Towards Riemannian accelerated gradient methods. *arXiv preprint arXiv:1806.02812*.