
Metric-Projected Accelerated Riemannian Optimization: Handling Constraints to Bound Geometric Penalties

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We propose an accelerated first-order method for the optimization of smooth and
2 (strongly or not) geodesically-convex functions over a compact and geodesically-
3 convex set in Hadamard manifolds, that we access to via a metric-projection oracle.
4 It enjoys the same rates of convergence as Nesterov’s accelerated gradient descent,
5 up to a multiplicative geometric penalty and log factors. Even without in-manifold
6 constraints, all prior fully accelerated works require their iterates to remain in
7 some specified compact set (which is needed in worst-case analyses due to a lower
8 bound), while only two previous methods are able to enforce this condition and
9 these, in contrast, have limited applicability, e.g., to local optimization or to spaces
10 of constant curvature. Our results solve an open question in [KY22] and an another
11 question related to one posed in [ZS16]. In our solution, we show we can use
12 projected Riemannian gradient descent to implement an inexact proximal point
13 operator that we use as a subroutine, which is of independent interest.

14 1 Introduction

15 Riemannian optimization concerns the optimization of a function defined over a Riemannian manifold.
16 It is motivated by constrained problems that can be naturally expressed on Riemannian manifolds
17 allowing to exploit the geometric structure of the problem and effectively transforming it into an
18 unconstrained one. Moreover, there are problems that are not convex in the Euclidean setting, but
19 that when posed as problems over a manifold with the right metric, are convex when restricted to
20 every geodesic, and this allows for fast optimization [Cru+06; CM12; BFO15; All+18]. That is, they
21 are geodesically convex (g-convex) problems, cf. Definition 1.1. Some applications of Riemannian
22 optimization in machine learning include low-rank matrix completion [CA16; HS18; MS14; Tan+14;
23 Van13], dictionary learning [CS17; SQW17], optimization under orthogonality constraints [EAS98;
24 LM19], robust covariance estimation in Gaussian distributions [Wie12], Gaussian mixture models
25 [HS15], operator scaling [All+18], and sparse principal component analysis [GHT15; HW19b;
26 JTU03].

27 Riemannian optimization, whether under g-convexity or not, is an extensive and active area of
28 research, for which one aspires to develop Riemannian optimization algorithms that share analogous
29 properties to the more broadly studied Euclidean first-order methods, such as the following kinds of
30 Riemannian methods: deterministic [BFM17; Wei+16; ZS16], adaptive [KJM19], projection-free
31 [WS17; WS19], saddle-point-escaping [CB19; SFF19; ZZS18; ZYF19; CB20], stochastic [HS17;

⁰Most of the notations in this work have a link to their definitions. For example, if you click or tap on any instance of L , you will jump to the place where it is defined as the smoothness constant of the function we consider in this work.

32 KL17; Tri+18], variance-reduced [SKM17; SKM19; ZRS16], and min-max methods [ZZS22], among
33 others.

34 Riemannian generalizations to accelerated convex optimization are appealing due to their better
35 convergence rates with respect to unaccelerated methods, specially in ill-conditioned problems.
36 Acceleration in Euclidean convex optimization is a concept that has been broadly explored and has
37 provided many different fast algorithms. A paradigmatic example is Nesterov’s Accelerated Gradient
38 Descent (AGD), cf. [Nes83], which can be considered the first general accelerated method, where
39 the conjugate gradients method can be seen as an accelerated predecessor in a more limited scope
40 [Mar21]. There have been recent efforts to better understand this phenomenon in the Euclidean case
41 [AO17; SBC16; DT14; WWJ16; DO19; Jou+20], which have yielded some fruitful techniques for
42 the general development of methods and analyses. These techniques have allowed for a considerable
43 number of new results going beyond the standard oracle model, convexity, or beyond first-order, in
44 a wide variety of settings [Tse08; BT09; WRM16; AO15; All17; All+16; All18b; Car+17; DO18;
45 All18a; CDO18; HSS19; CS19; DJ19; Gas+19; Iva+21; DN20; KG20; CMP21], among many others.
46 There have been some efforts to achieve acceleration for Riemannian algorithms as generalizations of
47 AGD, cf. Section 1.3. These works try to answer the following fundamental question:

48 *Can a Riemannian first-order method enjoy the same rates of convergence as Euclidean AGD?*

49 The question is posed under (possibly strongly) geodesic convexity and smoothness of the function to
50 be optimized. And we now know, due to the lower bound in [CB21], that the optimization should be
51 over a bounded domain and under bounded geodesic curvature of the Riemannian manifold. In this
52 work, we study this question in the case of Hadamard manifolds \mathcal{M} of bounded sectional curvature,
53 where many of the applications lie [HS20]. Given a compact and uniquely geodesic g -convex set \mathcal{X}
54 that we access to via a metric-projection oracle, we design first-order algorithms that enjoy the same
55 rates as AGD when approximating $\min_{x \in \mathcal{X}} f(x)$, up to logarithmic factors and up to a geometric
56 penalty factor, where $f : \mathcal{N} \subset \mathcal{M} \rightarrow \mathbb{R}$ is a differentiable function that is smooth and g -convex (or
57 strongly g -convex) in $\mathcal{X} \subset \mathcal{N}$. See Section 1.1 for the definitions of these concepts. Importantly,
58 our algorithm obtains acceleration without an undesirable assumption that most previous works
59 had to make: that the iterates of the algorithm stay inside of a specified compact set without any
60 mechanism for enforcing this condition. Only two previous methods are able to deal with some form
61 of constraints, and they apply to the limited settings of constant sectional curvature manifolds and
62 local optimization, respectively. Techniques in the rest of papers can handle neither constraints nor
63 projections, due to fundamental properties of their methods. Removing this condition in general,
64 global, and fully accelerated methods was posed as an open question in [KY22], that we solve for the
65 case of Hadamard manifolds. The difficulty of constraining problems in order to bound geometric
66 penalties as well as the necessity of achieving this goal in order to provide full optimization guarantees
67 is something that has also been noted in other kinds of Riemannian algorithms, cf. [HS20]. See
68 Table 1 for a succinct comparison among algorithms with some degree of acceleration and their rates.

69 The question concerning whether there are Riemannian analogs to Nesterov’s algorithm that enjoy
70 similar rates is a question that, to the best of our knowledge, was first formulated in [ZS16]. In
71 particular, since Nesterov’s AGD uses a proximal operator of a function’s linearization, they ask
72 whether there is a Riemannian analog to this operation that could be used to obtain accelerated rates
73 in the Riemannian case. The natural candidate results in a non-convex problem which is not amenable
74 to optimization. While we do not take this course of action, we show that, instead, a proximal step
75 with respect to the *whole* function can be approximated efficiently in Hadamard manifolds and it
76 can be used along with an accelerated outer loop, when implemented and analyzed carefully, in the
77 spirit of other Euclidean algorithms like Catalyst [LMH17]. It relies on Riemannian gradient descent
78 (RGD) with projections, initialized at a suitable warm-start point that we can find by exploiting the
79 structure of the geometry and the metric projection. The Riemannian proximal point subroutine
80 we design is of independent interest. To the best of our knowledge, previously known Riemannian
81 proximal methods either obtain asymptotic analyses, assume exact proximal computation, or work
82 with approximate proximal operators by using different inexactness conditions as ours, and do not
83 show how to implement the inexact operators, cf. Section 1.3.

84 1.1 Preliminaries

85 We provide definitions of Riemannian geometry concepts that we use in this work. The interested
86 reader can refer to [Pet06; Bac14] for an in-depth review of this topic, but for this work the following

87 notions will be enough. A Riemannian manifold $(\mathcal{M}, \mathfrak{g})$ is a real C^∞ manifold \mathcal{M} equipped with
88 a metric \mathfrak{g} , which is a smoothly varying, i.e., C^∞ , inner product. For $x \in \mathcal{M}$, denote by $T_x\mathcal{M}$ the
89 tangent space of \mathcal{M} at x . For vectors $v, w \in T_x\mathcal{M}$, we denote the inner product of the metric by
90 $\langle v, w \rangle_x$ and the norm it induces by $\|v\|_x \stackrel{\text{def}}{=} \sqrt{\langle v, v \rangle_x}$. Most of the time, the point x is known from
91 context, in which case we write $\langle v, w \rangle$ or $\|v\|$.

92 A geodesic of length ℓ is a curve $\gamma : [0, \ell] \rightarrow \mathcal{M}$ of unit speed that is locally distance minimizing.
93 A uniquely geodesic space is a space such that for every two points there is one and only one
94 geodesic that joins them. In such a case the exponential map $\text{Exp}_x : T_x\mathcal{M} \rightarrow \mathcal{M}$ and the inverse
95 exponential map $\text{Log}_x : \mathcal{M} \rightarrow T_x\mathcal{M}$ are well defined for every pair of points, and are as follows.
96 Given $x, y \in \mathcal{M}$, $v \in T_x\mathcal{M}$, and a geodesic γ of length $\|v\|$ such that $\gamma(0) = x$, $\gamma(\|v\|) = y$,
97 $\gamma'(0) = v/\|v\|$, we have that $\text{Exp}_x(v) = y$ and $\text{Log}_x(y) = v$. We denote by $d(x, y)$ the distance
98 between x and y , and note that it takes the same value as $\|\text{Log}_x(y)\|$. The manifold \mathcal{M} comes with a
99 natural parallel transport between vectors in different tangent spaces, that formally is defined from a
100 way of identifying nearby tangent spaces, known as the Levi-Civita connection ∇ [Lev77]. We use
101 this parallel transport throughout this work.

102 Given a 2-dimensional subspace $V \subseteq T_x\mathcal{M}$ of the tangent space of a point x , the sectional curvature
103 at x with respect to V is defined as the Gauss curvature, for the surface $\text{Exp}_x(V)$ at x . The Gauss
104 curvature at a point x can be defined as the product of the maximum and minimum curvatures of
105 the curves resulting from intersecting the surface with planes that are normal to the surface at x . A
106 Hadamard manifold is a complete simply connected Riemannian manifold whose sectional curvature
107 is non-positive, like the hyperbolic space or the space of $n \times n$ symmetric positive definite matrices
108 with the metric $\langle X, Y \rangle_A \stackrel{\text{def}}{=} \text{Tr}(A^{-1}XA^{-1}Y)$ where X, Y are in the tangent space of A . Hadamard
109 manifolds are uniquely geodesic. Note that in a general manifold $\text{Exp}_x(\cdot)$ might not be defined for
110 each $v \in T_x\mathcal{M}$, but in a Hadamard manifold of dimension n , the exponential map at any point is a
111 global diffeomorphism between $T_x\mathcal{M} \cong \mathbb{R}^n$ and the manifold, and so the exponential map is defined
112 everywhere. We now proceed to define the main properties that will be assumed on our model for the
113 function to be minimized and on the feasible set \mathcal{X} .

114 **Definition 1.1 (Geodesic Convexity and Smoothness).** Let $f : \mathcal{N} \subset \mathcal{M} \rightarrow \mathbb{R}$ be a differentiable
115 function defined on an open set \mathcal{N} contained in a Riemannian manifold \mathcal{M} . Given $L \geq \mu > 0$, we
116 say that f is L -smooth in \mathcal{X} if for any two points $x, y \in \mathcal{X}$, f satisfies

$$f(y) \leq f(x) + \langle \nabla f(x), \text{Log}_x(y) \rangle + \frac{L}{2}d(x, y)^2.$$

117 Analogously, we say that f is μ -strongly g -convex in \mathcal{X} , if for any two points $x, y \in \mathcal{X}$, we have

$$f(y) \geq f(x) + \langle \nabla f(x), \text{Log}_x(y) \rangle + \frac{\mu}{2}d(x, y)^2.$$

118 If the previous inequality is satisfied with $\mu = 0$, we say the function is g -convex in \mathcal{X} .

119 **Definition 1.2 (Metric projection operator).** Let \mathcal{M} be a Hadamard manifold and let $\mathcal{X} \subset \mathcal{M}$ be
120 a closed g -convex subset of \mathcal{M} . A *metric projection operator* onto \mathcal{X} is a map $\mathcal{P}_{\mathcal{X}} : \mathcal{M} \rightarrow \mathcal{X}$
121 satisfying $d(x, \mathcal{P}_{\mathcal{X}}(x)) \leq d(x, y)$ for all $y \in \mathcal{X}$.

122 A consequence of the definition is that the projection is single valued and non-expansive, the latter
123 meaning $d(\mathcal{P}_{\mathcal{X}}(x), \mathcal{P}_{\mathcal{X}}(y)) \leq d(x, y)$, cf. [Bac14, Thm 2.1.12].

124 We present the following fact about the squared distance function, when one of the arguments is fixed.
125 The constants ζ_D, δ_D below appear everywhere in Riemannian optimization because, among other
126 things, Fact 1.3 yields Riemannian inequalities that are analogous to the equality in the Euclidean
127 cosine law of a triangle, cf. Corollary B.3, and these inequalities have wide applicability in the
128 analyses of Riemannian methods.

129 **Fact 1.3 (Local information of the squared distance).** Let \mathcal{M} be a Riemannian manifold of sectional
130 curvature bounded by $[\kappa_{\min}, \kappa_{\max}]$ that contains a uniquely g -convex set $\mathcal{X} \subset \mathcal{M}$ of diameter
131 $D < \infty$. Then, given $x, y \in \mathcal{X}$ we have the following for the function $\Phi_x : \mathcal{M} \rightarrow \mathbb{R}, y \mapsto \frac{1}{2}d(x, y)^2$:

$$\nabla \Phi_x(y) = -\text{Log}_y(x) \quad \text{and} \quad \delta_D \|v\|^2 \leq \text{Hess } \Phi_x(y)[v, v] \leq \zeta_D \|v\|^2,$$

132 where

$$\zeta_D \stackrel{\text{def}}{=} \begin{cases} D\sqrt{|\kappa_{\min}|} \coth(D\sqrt{|\kappa_{\min}|}) & \text{if } \kappa_{\min} \leq 0 \\ 1 & \text{if } \kappa_{\min} > 0 \end{cases},$$

133 and

$$\delta_D \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } \kappa_{\max} \leq 0 \\ D\sqrt{\kappa_{\max}} \coth(D\sqrt{\kappa_{\max}}) & \text{if } \kappa_{\max} > 0 \end{cases},$$

134 In particular, Φ_x is δ_D -strongly g -convex and ζ_D -smooth in \mathcal{X} . See [Lez20] for a proof.

135 1.2 Notation.

136 Let \mathcal{M} be a uniquely geodesic n -dimensional Riemannian manifold. Given points $x, y, z \in \mathcal{M}$,
 137 we abuse the notation and write y in non-ambiguous and well-defined contexts in which we should
 138 write $\text{Log}_x(y)$. For example, for $v \in T_x\mathcal{M}$ we have $\langle v, y - x \rangle = -\langle v, x - y \rangle = \langle v, \text{Log}_x(y) -$
 139 $\text{Log}_x(x) \rangle = \langle v, \text{Log}_x(y) \rangle$; $\|v - y\| = \|v - \text{Log}_x(y)\|$; $\|z - y\|_x = \|\text{Log}_x(z) - \text{Log}_x(y)\|$; and
 140 $\|y - x\|_x = \|\text{Log}_x(y)\| = d(y, x)$. We denote by \mathcal{X} a compact, uniquely geodesic g -convex set of
 141 diameter D contained in an open set $\mathcal{N} \subset \mathcal{M}$ and we use $I_{\mathcal{X}}$ for the indicator function of \mathcal{X} , which
 142 is 0 at points in \mathcal{X} and $+\infty$ otherwise. For a vector $v \in T_y\mathcal{M}$, we use $\Gamma_y^x(v) \in T_x\mathcal{M}$ to denote the
 143 parallel transport of v from $T_y\mathcal{M}$ to $T_x\mathcal{M}$ along the unique geodesic that connects y to x . We call
 144 $f : \mathcal{N} \subset \mathcal{M} \rightarrow \mathbb{R}$ a differentiable L -smooth g -convex function we want to optimize over \mathcal{X} . We use
 145 ε to denote the approximation accuracy parameter, $x_0 \in \mathcal{X}$ for the initial point of our algorithms, and
 146 $R_0 \stackrel{\text{def}}{=} d(x_0, x^*)$ for the initial distance to an arbitrary minimizer $x^* \in \arg \min_{x \in \mathcal{X}} f(x)$. The big
 147 O notation $\tilde{O}(\cdot)$ omits log factors and $O^*(\cdot)$ omits log factors except those with respect to LR_0^2/ε .
 148 The latter will be useful to describe the rates of convergence for the strongly g -convex case, by
 149 emphasizing that there is no extra dependence on ε . Note that in the setting of Hadamard manifolds,
 150 the bounds on the sectional curvature are $\kappa_{\min} \leq \kappa_{\max} \leq 0$. Hence for convenience, given that we
 151 optimize over \mathcal{X} , we define $\zeta \stackrel{\text{def}}{=} \zeta_D = D\sqrt{|\kappa_{\min}|} \coth(D\sqrt{|\kappa_{\min}|}) \geq 1$ and $\delta \stackrel{\text{def}}{=} 1$. If $v \in T_x\mathcal{M}$,
 152 we use $\Pi_{\tilde{B}(0, D)}(v) \in T_x\mathcal{M}$ for the projection of v onto the closed ball with center at 0 and radius D .

153 1.3 Our results and comparisons with related work

154 In this work, we optimize functions defined over Hadamard manifolds \mathcal{M} of finite dimension n
 155 and of sectional curvature bounded in $[\kappa_{\min}, \kappa_{\max}]$. As all previous related works discussed in the
 156 sequel, we assume that we can compute the exponential and inverse exponential maps, and parallel
 157 transport of vectors for our manifold. The differentiable function f to be optimized is defined over
 158 an open set $\mathcal{N} \subset \mathcal{M}$ that contains a compact g -convex set \mathcal{X} of finite diameter D , that we access
 159 via a metric-projection oracle. Our function f is L -smooth and g -convex (or μ -strongly g -convex)
 160 in \mathcal{X} and we have access to it via a gradient oracle that can be queried at points in \mathcal{X} . For the
 161 setting we just described, we show in [Theorem 2.2](#) and [Theorem 2.4](#) that the algorithms we propose
 162 find a point $y_T \in \mathcal{X}$ such that $f(y_T) - \min_{x \in \mathcal{X}} f(x) \leq \varepsilon$ after calling the gradient oracle and the
 163 metric-projection oracle the following number of times: $\tilde{O}(\zeta^2 \sqrt{LR_0^2/\varepsilon})$ for the g -convex case and
 164 $O^*(\zeta^2 \sqrt{L/\mu} \log(\mu R_0^2/\varepsilon))$ for the μ -strongly g -convex case, where $R_0 \stackrel{\text{def}}{=} d(x_0, x^*)$ and $x_0 \in \mathcal{X}$ is
 165 an initial point. That is, the algorithms enjoy the same rates as AGD in the Euclidean space up to a
 166 factor of $\zeta^2 = D^2 \kappa_{\min}^2 \coth^2(D\sqrt{|\kappa_{\min}|})$ (our geometric penalty) and up to universal constants and
 167 log factors. Note that as the minimum curvature κ_{\min} approaches 0 we have $\zeta \rightarrow 1$.

168 We emphasize that our algorithms only need to query the gradient of f at points in \mathcal{X} and the
 169 L -smoothness and μ -strong g -convexity of f only need to hold in \mathcal{X} . This is relevant because in
 170 Riemannian manifolds the condition number L/μ in a set can increase with the size of the set, cf.
 171 [\[Mar22, Proposition 27\]](#). Intuitively, although there are twice differentiable functions defined over the
 172 Euclidean space whose Hessian is constant everywhere, in other Riemannian cases the metric may
 173 preclude having such global condition and the larger the (compact) set is, the greater the maximum
 174 eigenvalue of the Hessian over this set (i.e., its smoothness constant) can be with respect to the
 175 minimum one (i.e., its strongly g -convex constant) for any smooth and strongly g -convex function.
 176 Compare this, for instance, with the bounds on the Hessian's eigenvalues of the squared distance
 177 function in [Fact 1.3](#), which are tight for spaces of constant curvature [\[Lez20\]](#).

178 Now we proceed to compare our results with previous works. We have summarized most of the
 179 following discussion in [Table 1](#). We include Nesterov's AGD in the table for comparison purposes¹.

¹Note that the original method in [\[Nes83\]](#) needed to query the gradient of the function outside of the feasible set, and this was later improved to only require queries at feasible points [\[Nes05\]](#) as in our work, hence our choice of citation in the table.

180 There are some works on Riemannian acceleration that focus on empirical evaluation or that work
181 under strong assumptions [Liu+17; Ali+19; HW19a; Ali+20; Lin+20], see [Mar22] for instance for
182 a discussion on these works. We focus the discussion on the most related work with guarantees.
183 [ZS18] obtain an algorithm that, up to constants, achieves the same rates as AGD in the Euclidean
184 space, for L -smooth and μ -strongly g -convex functions but only *locally*, namely when the initial
185 point starts in a small neighborhood N of the minimizer x^* : a ball of radius $O((\mu/L)^{3/4})$ around it.
186 [AS20] generalize the previous algorithm and, by using similar ideas as in [ZS18] for estimating a
187 lower bound on f , they adapt the algorithm to work globally, proving that it eventually decreases the
188 objective as fast as AGD. However, as [Mar22] noted, it takes as many iterations as the ones needed
189 by RGD to reach the neighborhood of the previous algorithm. The latter work also noted that in fact
190 RGD and the algorithm in [ZS18] can be run in parallel and combined to obtain the same convergence
191 rates as in [AS20], which suggested that for this technique, full acceleration with the rates of AGD
192 only happens over the small neighborhood N in [ZS18]. Note however that [AS20] show that
193 their algorithm will decrease the function value faster than RGD, but this is not quantified. [JS21]
194 developed a different framework, arising from [AS20] but with the same guarantees for accelerated
195 first-order methods. We do not feature it in the table. [CB21] showed that in a ball of center $x \in \mathcal{M}$
196 and radius $O((\mu/L)^{1/2})$ containing x^* , the pullback function $f \circ \text{Exp}_x : T_x \mathcal{M} \rightarrow \mathbb{R}$ is strongly
197 convex and smooth with condition number $O(L/\mu)$, so they argue that using AGD on the pullback
198 over the corresponding pulled-back Euclidean ball in the tangent space results in local acceleration
199 as well. In short, acceleration is possible in a small neighborhood because there the manifold is
200 almost Euclidean and the geometric deformations are small in comparison to the curvature of the
201 objective. These techniques do not work with the g -convex case since the neighborhood becomes a
202 point ($\mu/L = 0$).

203 Finding fully accelerated algorithms that are *global* presents a harder challenge. By a fully accelerated
204 algorithm we mean one with rates with same dependence as AGD on L , ε , and if it applies, on μ .
205 [Mar22] provided such algorithms for g -convex functions, strongly or not, defined over manifolds of
206 constant sectional curvature and constrained to a ball of radius R . In the convergence rates, there is a
207 geometric factor of $c = \cos(R\sqrt{K})^{-\Theta(1)}$ for sectional curvature $K > 0$, and $c = \cosh(R\sqrt{-K})^{\Theta(1)}$
208 when $K < 0$, cf. Table 1. When $R\sqrt{|K|} = O(1)$, they recover the same rates as AGD, which for
209 those manifolds is more general than the local assumption in the previous set of works. For larger
210 values of $R\sqrt{|K|}$, there is also full acceleration, but note that c grows rapidly when $K < 0$, since
211 there is an exponential dependence on R . When $K > 0$ the geometric penalty also grows fast, but
212 this is more natural since the minimum condition number of a function in a ball of radius R grows
213 similarly [Mar22]. The geometric penalties are large in some regimes because the algorithm bounds
214 uniformly, over the whole domain, the worst-case deformations that can occur. On the other hand, for
215 manifolds of bounded sectional curvature, [KY22] design algorithms with the same rates as AGD
216 up to universal constants and a factor of ζ , their geometric penalty. However, they need to assume
217 that the iterates of their algorithm remain in \mathcal{X} and point out on the necessity of removing such an
218 assumption, which they leave as an open question. Our work solves this question for the case of
219 Hadamard manifolds. In their technique, they show that they can use the structure of the accelerated
220 scheme to *move* lower bound estimations on $f(x^*)$ from one particular tangent space to another
221 without incurring extra errors, when the right Lyapunov function is used. By *moving* lower bounds
222 here we mean finding suitable lower bounds that are simple (a quadratic in their case), if pulled-back
223 to one tangent space, if we start with a similar bound that is simple when pulled-back to another
224 tangent space.

225 **Lower bounds.** In this paragraph, we omit constants depending on the curvature bounds in the
226 big- O notations for simplicity. [HM21] proved an optimization lower bound showing that acceleration
227 in Riemannian manifolds is harder than in the Euclidean space. [CB21] largely generalized their
228 results. They essentially show that for a large family of Hadamard manifolds, there is a function
229 that is smooth and strongly g -convex in a ball of radius R that contains the minimizer x^* , and for
230 which finding a point that is $R/5$ close to x^* requires $\tilde{\Omega}(R)$ calls to the gradient oracle. Note that
231 these results do not preclude the existence of a fully accelerated algorithm with rates $\tilde{O}(R)$ +AGD
232 rates, for instance. But they show that even if we want to perform unconstrained optimization, so
233 no in-manifold constraints are originally imposed, we need to optimize over a bounded domain in
234 order to bound geometric penalties. A similar statement is provided in the case of smooth and only
235 g -convex functions.

Table 1: Convergence rates of related works with provable guarantees for smooth problems over uniquely geodesic manifolds, in chronological order with respect to when the works were publicly available. Column **K?** refers to the supported values of the sectional curvature, **G?** to whether the algorithm is global (any initial distance to a minimizer is allowed). Here L and L' mean they are local algorithms that require initial distance $O((L/\mu)^{-3/4})$ and $O((L/\mu)^{-1/2})$, respectively. Column **F?** refers to whether there is full acceleration, meaning dependence on L , μ , and ε like AGD up to possibly log factors. Column **C?** refers to whether the method supports constraints. All methods require their iterates to be in some specified compact set, but the works with \times just assume the iterates will remain within the constraints, while the ones with \checkmark can force this condition with a projection oracle. Also, here **B** is like \checkmark but with the constraints limited to a ball. See Section 1.3 for the value c in [Mar22]. We use $\mathcal{W} \stackrel{\text{def}}{=} \sqrt{\frac{L}{\mu}} \log\left(\frac{LR_0^2}{\varepsilon}\right)$. *In [CB21], a condition is required on the covariant derivative of the metric tensor, cf. [CB21, Section 6].

| Method | g-convex | μ -st. g-convex | K? | G? | F? | C? |
|-------------------------------|---|------------------------------------|-----------------|--------------|--------------|--------------|
| [Nes05, AGD] | $O(\sqrt{\frac{LR_0^2}{\varepsilon}})$ | $O(\mathcal{W})$ | 0 | \checkmark | \checkmark | \checkmark |
| [ZS18, Theorem 11] | - | $O(\mathcal{W})$ | bounded | L | \checkmark | \times |
| [AS20, Theorem 3.1] | - | $O^*(\frac{L}{\mu} + \mathcal{W})$ | bounded | \checkmark | \times | \times |
| [Mar22, Remark 30] | - | $O^*(\frac{\mu}{L} + \mathcal{W})$ | bounded | \checkmark | \times | \times |
| [Mar22, Theorems 6 & 8] | $\tilde{O}(c\sqrt{\frac{LR_0^2}{\varepsilon}})$ | $O^*(c \cdot \mathcal{W})$ | ctant. $\neq 0$ | \checkmark | \checkmark | B |
| [CB21, Section 6] | - | $O(\mathcal{W})$ | bounded* | L' | \checkmark | B |
| [KY22, Corollaries 1 & 2] | $O(\zeta\sqrt{\frac{LR_0^2}{\varepsilon}})$ | $O(\zeta \cdot \mathcal{W})$ | bounded | \checkmark | \checkmark | \times |
| Theorems 2.2 & 2.4 | $\tilde{O}(\zeta^2\sqrt{\frac{LR_0^2}{\varepsilon}})$ | $O^*(\zeta^2 \cdot \mathcal{W})$ | Hadamard | \checkmark | \checkmark | \checkmark |

236 **Handling constraints to bound geometric penalties.** Due to the lower bounds, it becomes crucial
237 for a fully accelerated algorithm to restrict the optimization to a set \mathcal{X} of finite diameter D , or
238 otherwise a worst-case analysis incurs an arbitrary large geometric penalty in the rates. In our
239 algorithm and in all other known fully accelerated algorithms, learning rates depend on this diameter.
240 This is natural: estimation errors due to geometric deformations depend on the diameter via the
241 constants ζ_D , δ_D , the cosine-law inequalities Corollary B.3, or other analogous inequalities, and the
242 algorithms take these errors into account. All other previous works are not able to deal with any
243 constraints and hence they simply assume that the iterates of their algorithms stay within one such
244 specified set, except for [Mar22] and [CB21] that enforce a ball constraint, as we explained above.
245 However, these two works have their applicability limited to spaces of constant curvature and to local
246 optimization, respectively. Note that even if one could show in some settings that given a choice of
247 learning rate, convergence implies that the iterates will remain in some compact set, then because
248 the learning rates depend on the diameter of the set, and the diameter of the set would depend on
249 the learning rates, one cannot conclude from this argument that the assumption these works make is
250 going to be satisfied. In contrast, in this work, we design the first accelerated algorithm that supports
251 metric projections and, consequently, we can handle general constraints to bound geometric penalties
252 and accelerate our method without any other extra assumptions, solving an open question in [KY22].

253 Some other works study and use Riemannian metric projections in other contexts, see [Wal74;
254 HP13; BHP13; Bac14; ZS16] and references therein. Among them, [ZS16] introduced several, both
255 deterministic and stochastic, *unaccelerated* first-order methods that work with in-manifold constraints
256 by using metric-projection oracles. Our Algorithm 1 uses their projected RGD as a subroutine, cf.
257 Remark 2.3.

258 **Finding a global minimizer.** In our work, we do not need to assume that the set \mathcal{X} contains a
259 global minimizer, namely a point x^* such that $\nabla f(x^*) = 0$. We find an ε -minimizer with respect to
260 the minimum value of f at \mathcal{X} . All other previous works assume that the set contains the minimizers
261 of f , with the exception of [Mar22], where the algorithm can forgo this assumption if one has
262 access to a bound $L_{f,\mathcal{B}}$ on the Lipschitz constant of f when restricted to their ball constraint \mathcal{B} , and
263 in such a case the rates have a $\log(L_{f,\mathcal{B}}D/\varepsilon)$ factor instead of a $\log(LD^2/\varepsilon)$ factor. Note this is
264 natural since if a global minimizer is in the set, then we have $L_{f,\mathcal{B}} = O(LD)$. We note that we also

265 obtain a logarithmic dependence that involves the Lipschitz constant $L_{f,\mathcal{X}}$ of f in \mathcal{X} (the logarithmic
 266 dependence involves the scale invariant quantity ζ_C for $C = L_{f,\mathcal{X}}/L$, which is $O(\zeta)$ if $x^* \in \mathcal{X}$) but
 267 in contrast in our case, our method does not require access to the Lipschitz constant of f in \mathcal{X} .

268 **Riemannian proximal methods** There have been some works that study proximal methods in
 269 Riemannian manifolds, but most of them focus on asymptotic results or assume the proximal operator
 270 can be computed exactly [Wan+15; BFM17; BCO16; Kha+21; Cha+21]. The rest of these works
 271 study proximal point methods under different inexact versions of the proximal operator as ours and
 272 they do not show how to implement their inexact version in applications, like our case of smooth and
 273 g-convex optimization. [AK14] provide a convergence analysis of an inexact proximal point method
 274 but when applied to optimization they assume the computation of the proximal operator is exact.
 275 [TH14] uses a different inexact condition and proves linear convergence, under a growth condition
 276 on f . [Wan+16] obtains linear convergence of an inexact proximal point method under a different
 277 growth assumption on f and under an absolute error condition on the proximal function.

278 2 Algorithm and Pseudocode

279 In this section, we present our **Riemannian accelerated algorithm** for **constrained g-convex** optimiza-
 280 tion, or **Riemacon**². Recall our abuse of notation for points $p \in \mathcal{M}$ to mean $\text{Log}_q(p)$ in contexts in
 281 which one should place a vector in $T_q\mathcal{M}$ and note that in our algorithm x_k and y_k are points in \mathcal{M}
 282 whereas $z_k^{x_k} \in T_{x_k}\mathcal{M}$, $z_k^{y_k}, \bar{z}_k^{y_k} \in T_{y_k}\mathcal{M}$.

Algorithm 1 Riemacon: Riemannian Acceleration - Constrained g-Convex Optimization

Input: Initial point $x_0 \in \mathcal{X} \subset \mathcal{N}$. Diff. function $f : \mathcal{N} \subset \mathcal{M} \rightarrow \mathbb{R}$ for a Hadamard manifold \mathcal{M}
 that is L -smooth and g-convex in \mathcal{X} , final iteration T (not required to be known in advance).

Parameters:

- Geometric penalty $\xi \stackrel{\text{def}}{=} 4\zeta_{2D} - 3 \leq 8\zeta - 3 = O(\zeta)$.
- Implicit Gradient Descent learning rate $\lambda \stackrel{\text{def}}{=} \zeta_{2D}/L$.
- Mirror Descent learning rates $\eta_k \stackrel{\text{def}}{=} a_k/\xi$.
- Proportionality constant in the proximal subproblem accuracies: $\Delta_k \stackrel{\text{def}}{=} \frac{1}{(k+1)^2}$.

Definition: (computation of this value is not needed)

- Prox. accuracies: $\sigma_k \stackrel{\text{def}}{=} \frac{\Delta_k d(x_k, y_k^*)^2}{78\lambda}$ where $y_k^* \stackrel{\text{def}}{=} \arg \min_{y \in \mathcal{X}} \{f(y) + \frac{1}{2\lambda} d(x_k, y)^2\}$.
-

```

1:  $y_0 \leftarrow x_0$ ;  $A_0 \leftarrow 200\lambda\xi$ 
2:  $z_0^{x_0} \leftarrow 0 \in T_{x_0}\mathcal{M}$ ;  $\bar{z}_0^{y_0} \leftarrow z_0^{y_0} \leftarrow 0 \in T_{y_0}\mathcal{M}$ 
3: for  $k = 1$  to  $T$  do
4:    $a_k \leftarrow 2\lambda \frac{k+32\xi}{5}$ 
5:    $A_k \leftarrow a_k/\xi + A_{k-1} = \sum_{i=1}^k a_i/\xi + A_0 = \lambda \left( \frac{k(k+1+64\xi)}{5\xi} + 200\xi \right)$ 
6:    $x_k \leftarrow \text{Exp}_{y_{k-1}} \left( \frac{a_k}{A_{k-1}+a_k} \bar{z}_{k-1}^{y_{k-1}} + \frac{A_{k-1}}{A_{k-1}+a_k} y_{k-1} \right) = \text{Exp}_{y_{k-1}} \left( \frac{a_k}{A_{k-1}+a_k} \bar{z}_{k-1}^{y_{k-1}} \right)$   $\diamond$  Coupling
7:    $z_{k-1}^{x_k} \leftarrow \Gamma_{y_{k-1}}^{x_k}(\bar{z}_{k-1}^{y_{k-1}}) + \text{Log}_{x_k}(y_{k-1}) = \text{Log}_{x_k}(\text{Exp}_{y_k}(\bar{z}_{k-1}^{y_{k-1}}))$ 
8:    $y_k \leftarrow \sigma_k$ -minimizer of the proximal problem  $\min_{y \in \mathcal{X}} \{f(y) + \frac{1}{2\lambda} d(x_k, y)^2\}$  (cf. Remark 2.3).
9:    $v_k^x \leftarrow -\text{Log}_{x_k}(y_k)/\lambda$   $\diamond$  Approximate subgradient
10:   $z_k^{x_k} \leftarrow z_{k-1}^{x_k} - \eta_k v_k^x$   $\diamond$  Mirror Descent step
11:   $z_k^{y_k} \leftarrow \Gamma_{x_k}^{y_k}(z_k^{x_k}) + \text{Log}_{y_k}(x_k)$   $\diamond$  Moving the dual point to  $T_{y_k}\mathcal{M}$ 
12:   $\bar{z}_k^{y_k} \leftarrow \Pi_{\bar{B}(0,D)}(z_k^{y_k}) \in T_{y_k}\mathcal{M}$   $\diamond$  Easy projection done so the dual point is not very far
13: end for
14: return  $y_T$ .
```

283 We start with an interpretation of our algorithm that helps understanding its high-level ideas. The fol-
 284 lowing intends to be a qualitative explanation, and we refer to the pseudocode and the supplementary
 285 material for the exact descriptions and analysis. Euclidean accelerated algorithms can be interpreted,
 286 cf. [AO17], as a combination of a gradient descent (GD) algorithm and an online learning algorithm

²Riemacon rhymes with “rima con” in Spanish.

287 with losses being the affine lower bounds $f(x_k) + \langle \nabla f(x_k), \cdot - x_k \rangle$ we obtain on $f(\cdot)$ by applying
288 convexity at some points x_k . That is, the latter builds a lower bound estimation on f . By selecting
289 the next query to the gradient oracle as a cleverly picked convex combination of the predictions given
290 by these two algorithms, one can show that the instantaneous regret of the online learning algorithm
291 can be compensated by the local progress GD makes, which leads to accelerated convergence. In
292 Riemannian optimization, there are two main obstacles. Firstly, the first-order approximations of f
293 at points x_k yield functions that are affine but only with respect to their respective $T_{x_k} \mathcal{M}$, and so
294 combining these lower bounds that are only simple in their tangent spaces makes obtaining good
295 global estimations not simple. Secondly, when one obtains such global estimations, then one naturally
296 incurs an instantaneous regret that is worse by a factor than is usual in Euclidean acceleration. This
297 factor is a geometric constant depending on the diameter D of a set \mathcal{X} where the iterates and a
298 (possibly constrained) minimizer lie. As a consequence, the learning rate of GD would need to be
299 multiplicatively increased by such a constant with respect to the one of the online learning algorithm
300 in order for the regret to still be compensated with the local progress of GD (and the rates worsen by
301 this constant). But if we fix some \mathcal{X} of finite diameter, because GD's learning rate is now larger, it is
302 not clear how to keep the iterates in \mathcal{X} . And if we do not have the iterates in one such set \mathcal{X} , then our
303 geometric penalties could grow arbitrarily.

304 We find the answer in implicit methods. An implicit Euclidean (sub)gradient descent step is one that
305 computes, from a point $x_k \in \mathcal{X}$, another point $y_k^* = x_k - \lambda v_k \in \mathcal{X}$, where $v_k \in \partial(f + I_{\mathcal{X}})(y_k^*)$, is a
306 subgradient of $f + I_{\mathcal{X}}$ at y_k^* . Intuitively, if we could implement a Riemannian version of an implicit
307 GD step then it should be possible to still compensate the regret of the other algorithm and keep all the
308 iterates in the set \mathcal{X} . Computing such an implicit step is computationally hard in general, but we show
309 that approximating the proximal objective $h_k(y) \stackrel{\text{def}}{=} f(y) + \frac{1}{2\lambda} d(x_k, y)^2$ with enough accuracy yields
310 an approximate subgradient that can be used to obtain an accelerated algorithm as well. In particular,
311 we provide an accelerated scheme for which we show that the error incurred by the approximation
312 of the subgradient can be bounded by some terms we can control, cf. [Lemma A.2](#), namely a small
313 term that appears in our Lyapunov function and also a term proportional to the squared norm of
314 the approximated subgradient, which only adds a constant to the final convergence rates. We also
315 provide a warm start in [Lemma A.4](#) and an analysis that shows that using the projected Riemannian
316 gradient descent in [\[ZS18\]](#) initialized at the warm-started point achieves the desired accuracy of
317 the subproblem fast, cf. [Remark 2.3](#). This proximal approach works by exploiting the fact that the
318 Riemannian Moreau envelop is convex in Hadamard manifolds [\[AF05\]](#) and that the subproblem h_k ,
319 defined with our $\lambda = \zeta_{2D}/L$, is strongly g -convex and smooth with a condition number that only
320 depends on the geometry. Besides of these steps, we use a coupling of the approximate implicit RGD
321 and of a mirror descent (MD) algorithm, along with a technique in [\[KY22\]](#) to move dual points to
322 the right tangent spaces without incurring extra geometric penalties, that we adapt to work with dual
323 projections, cf. [Lemma A.3](#). Importantly, the MD algorithm keeps the dual point close to the set \mathcal{X} by
324 using the projection in [Line 12](#), which implies that the point x_k is close to \mathcal{X} as well, and this is crucial
325 to keep low geometric penalties. This MD approach is a mix between follow-the-regularized-leader
326 algorithms, that do not project the dual variable, and pure mirror descent algorithms that always
327 project the dual variable. In the analysis, we note that partial projection also works, meaning that
328 defining a new dual point that is closer to all of the points in the feasible set but without being a full
329 projection leads to the same guarantees. Because we use the mirror descent lemma over $T_{y_k} \mathcal{M}$, what
330 we described translates to: we can project the dual $z_k^{y_k}$ onto a ball defined on $T_{y_k} \mathcal{M}$ that contains the
331 pulled-back set $\text{Log}_{y_k}(\mathcal{X})$ and by means of that trick we can keep the iterates x_k close to \mathcal{X} . And at
332 the same time, the point for which we prove guarantees, namely y_k , is always in \mathcal{X} .

333 We leave the proofs of most of our results to the supplementary material and state our main theorems
334 below. Using the insights explained above, we show the following inequality on ψ_k , defined below,
335 that will be used as a Lyapunov function to prove the convergence rates of [Algorithm 1](#).

336 **Proposition 2.1.** [\downarrow] *By using the notation of [Algorithm 1](#), let*

$$\psi_k \stackrel{\text{def}}{=} A_k(f(y_k) - f(x^*)) + \frac{1}{2} \|z_k^{y_k} - x^*\|_{y_k}^2 + \frac{\xi - 1}{2} \|y_k - z_k^{y_k}\|_{y_k}^2.$$

337 *Then, for all $k \geq 1$, we have $(1 - \Delta_k)\psi_k \leq \psi_{k-1}$.*

338 Finally, we can state our theorem for the optimization of L -smooth and g -convex functions.

339 **Theorem 2.2.** [\downarrow] *Let \mathcal{M} be a finite-dimensional Hadamard manifold of bounded sectional curvature,*
340 *let $f : \mathcal{N} \subset \mathcal{M} \rightarrow \mathbb{R}$ be an L -smooth and g -convex differentiable function in a compact g -convex*

341 set $\mathcal{X} \subset \mathcal{N}$ of diameter D , and $x^* \in \arg \min_{x \in \mathcal{X}} f(x)$. For $R_0 \stackrel{\text{def}}{=} d(x_0, x^*)$, and all $k \geq 1$,
 342 the iterates y_k of [Algorithm 1](#) satisfy $y_k \in \mathcal{X}$ and $f(y_k) - f(x^*) = O\left(\frac{LR_0^2}{k^2} \cdot \zeta^2\right)$. That is, after
 343 $T = O(\zeta \sqrt{\frac{LR_0^2}{\varepsilon}})$ iterations we find an ε -minimizer. Moreover, the total number of queries to the
 344 gradient and projection oracles is bounded by $\tilde{O}(\zeta^2 \sqrt{\frac{LR_0^2}{\varepsilon}})$.

345 We note that a straightforward corollary from our results is that if we can compute the exact Rie-
 346 mannian proximal point operator and we use it as the implicit gradient descent step in [Line 8](#) of
 347 [Algorithm 1](#), then the method is an accelerated proximal point method. One such Riemannian
 348 algorithm was previously unknown in the literature as well.

349 Now we show that [Line 8](#) can be implemented efficiently. The essential part is being able to have and
 350 use a point with the guarantees of our warm start, cf. [Lemma A.4](#).

351 **Remark 2.3 (Solving the subproblems).** Let \mathcal{A} be the unaccelerated Riemannian gradient descent
 352 algorithm in [[ZS16](#), [Theorem 15](#)]. This algorithm takes a function $h : \mathcal{M} \rightarrow \mathbb{R}$ with minimizer at y^*
 353 when restricted to $\mathcal{X} \subset \mathcal{M}$ that is μ' -strongly g -convex and L' -smooth in \mathcal{X} , where \mathcal{M} is a Hadamard
 354 manifold of bounded sectional curvature and \mathcal{X} is a geodesically-convex compact set with diameter
 355 D and returns a point p_t satisfying $h_k(p_t) - h_k(y^*) \leq \varepsilon'$ after querying a gradient oracle for h_k and
 356 a metric-projection oracle $\mathcal{P}_{\mathcal{X}}$ for \mathcal{X} for $t = O\left(\left(\zeta + \frac{L'}{\mu'}\right) \log\left(\frac{(h_k(p_0) - h_k(y^*)) + L'd(p_0, y^*)^2}{\varepsilon'}\right)\right)$ times³.

357 If we apply this algorithm to $h \leftarrow h_k(y) \stackrel{\text{def}}{=} f(y) + \frac{1}{2\lambda}d(x_k, y)^2$, we have $y^* \leftarrow y_k^*$, $L' \leftarrow 2L$ and
 358 $\mu' \leftarrow L/\zeta_{2D}$, so the condition number is $L'/\mu' = O(\zeta_{2D}) = O(\zeta)$. This is computed taking into
 359 account that f is L -smooth and 0 -strongly g -convex and using the ζ_{2D}/λ -smoothness and $1/\lambda$ -strong
 360 g -convexity of the second summand, which is given by [Fact 1.3](#) and [\(1\)](#). If we initialize the method with
 361 $p_0 \stackrel{\text{def}}{=} \mathcal{P}_{\mathcal{X}}(\text{Exp}_{x'_k}(-\frac{1}{L'}\nabla h_k(x'_k)))$, where $x'_k \stackrel{\text{def}}{=} \mathcal{P}_{\mathcal{X}}(x_k)$, then using (L/ζ_{2D}) -strong g -convexity of
 362 h_k to bound $L'd(p_0, y_k^*)^2 \leq 4\zeta_{2D}(h_k(p_0) - h(y_k^*))$, using [Lemma A.4](#) with $x \leftarrow x_k$, $p \leftarrow y_k^*$, and
 363 using the guarantees on \mathcal{A} , we have that we find a point y_k satisfying $h_k(y_k) - h_k(y_k^*) \leq \frac{\Delta_k d(x_k, y_k^*)^2}{78\lambda}$
 364 in $\tilde{O}(\zeta)$ queries to the gradient and projection oracles. See [Remark A.5](#) for the computation of this
 365 value. We note that any other algorithm with linear convergence rates for constrained strongly
 366 g -convex, smooth problems that works with a metric-projection oracle can be used as a subroutine to
 367 obtain an accelerated Riemannian algorithm.

368 We introduce the algorithm for μ -strongly g -convex functions via a reduction to [Algorithm 1](#), for
 369 simplicity. We note that the reverse Riemannian reduction yields extra factors in the rates depending
 370 on R_0 and the curvature, but this reduction does not yield any extra factors in the rates and in fact,
 371 it is slightly better than the usual convergence that is obtained when one analyzes these kinds of
 372 accelerated algorithms directly, by having a μ factor instead of L inside of the logarithm.

373 **Theorem 2.4.** [[↓](#)] Under the same assumptions as in [Theorem 2.2](#), let now f be μ -strongly g -convex.
 374 Applying the reduction in [[Mar22](#), [Theorem 7](#)], we obtain an algorithm that finds an ε -minimizer of f
 375 by querying the gradient oracle and projection oracle $O^*\left(\zeta^2 \sqrt{\frac{L}{\mu}} \log\left(\frac{\mu R_0^2}{\varepsilon}\right)\right)$ times.

376 3 Conclusion and future directions

377 In this work, we pursued an approach that, by designing inexact Riemannian proximal methods,
 378 yielded accelerated optimization algorithms that can work with metric projection oracles. Conse-
 379 quently we were able to work without an undesirable assumption that most previous methods required,
 380 whose potential satisfiability is not clear: that the iterates stay in certain specified geodesically-convex
 381 set without enforcing them to be in the set. A future direction of research is the study of whether there
 382 are algorithms like ours that incur even lower geometric penalties or that do not incur $\log(1/\varepsilon)$ factors.
 383 Another interesting direction consists of studying generalizations of our approach to manifolds of
 384 non-negative or of bounded sectional curvature manifolds.

³In their theorem, the authors only stated that $O\left(\left(\zeta + \frac{L'}{\mu'}\right) \log\left(\frac{L'D^2}{\varepsilon'}\right)\right)$ queries to the gradient oracle are enough, but in their proof they show this more refined statement, that we use.

385 References

- 386 [AF05] Daniel Azagra and Juan Ferrera. “Inf-convolution and regularization of convex functions
387 on Riemannian manifolds of nonpositive curvature”. In: *arXiv preprint math/0505496*
388 (2005) (cit. on p. 8).
- 389 [AK14] P Ahmadi and H Khatibzadeh. “On the convergence of inexact proximal point algorithm
390 on Hadamard manifolds”. In: *Taiwanese Journal of Mathematics* 18.2 (2014), pp. 419–
391 433 (cit. on p. 7).
- 392 [Ali+19] Foivos Alimisis, Antonio Orvieto, Gary Bécigneul, and Aurelien Lucchi. “A Continuous-
393 time Perspective for Modeling Acceleration in Riemannian Optimization”. In: *arXiv*
394 *preprint arXiv:1910.10782* (2019) (cit. on p. 5).
- 395 [Ali+20] Foivos Alimisis, Antonio Orvieto, Gary Bécigneul, and Aurelien Lucchi. “Practical Ac-
396 celerated Optimization on Riemannian Manifolds”. In: *arXiv preprint arXiv:2002.04144*
397 (2020) (cit. on p. 5).
- 398 [All+16] Zeyuan Allen Zhu, Zheng Qu, Peter Richtárik, and Yang Yuan. “Even Faster Acceler-
399 ated Coordinate Descent Using Non-Uniform Sampling”. In: *Proceedings of the 33rd*
400 *International Conference on Machine Learning, ICML 2016, New York City, NY, USA,*
401 *June 19-24, 2016*. 2016, pp. 1110–1119 (cit. on p. 2).
- 402 [All+18] Zeyuan Allen-Zhu, Ankit Garg, Yuanzhi Li, Rafael Oliveira, and Avi Wigderson. “Op-
403 erator scaling via geodesically convex optimization, invariant theory and polynomial
404 identity testing”. In: *Proceedings of the 50th Annual ACM SIGACT Symposium on*
405 *Theory of Computing*. 2018, pp. 172–181 (cit. on p. 1).
- 406 [All17] Zeyuan Allen-Zhu. “Katyusha: The First Direct Acceleration of Stochastic Gradient
407 Methods”. In: *J. Mach. Learn. Res.* 18 (2017), 221:1–221:51 (cit. on p. 2).
- 408 [All18a] Zeyuan Allen-Zhu. “Katyusha X: Practical Momentum Method for Stochastic Sum-of-
409 Nonconvex Optimization”. In: *Proceedings of the 35th International Conference on*
410 *Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15,*
411 *2018*. 2018, pp. 179–185 (cit. on p. 2).
- 412 [All18b] Zeyuan Allen-Zhu. “Natasha 2: Faster Non-Convex Optimization Than SGD”. In:
413 *Advances in Neural Information Processing Systems 31: Annual Conference on Neural*
414 *Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal,*
415 *Canada*. 2018, pp. 2680–2691 (cit. on p. 2).
- 416 [AO15] Zeyuan Allen Zhu and Lorenzo Orecchia. “Nearly-Linear Time Positive LP Solver
417 with Faster Convergence Rate”. In: *Proceedings of the Forty-Seventh Annual ACM on*
418 *Symposium on Theory of Computing, STOC 2015, Portland, OR, USA, June 14-17, 2015.*
419 2015, pp. 229–236 (cit. on p. 2).
- 420 [AO17] Zeyuan Allen Zhu and Lorenzo Orecchia. “Linear Coupling: An Ultimate Unification
421 of Gradient and Mirror Descent”. In: *8th Innovations in Theoretical Computer Science*
422 *Conference, ITCS 2017, January 9-11, 2017, Berkeley, CA, USA*. 2017, 3:1–3:22 (cit. on
423 pp. 2, 7).
- 424 [AS20] Kwangjun Ahn and Suvrit Sra. “From Nesterov’s Estimate Sequence to Riemannian
425 Acceleration”. In: *arXiv preprint arXiv:2001.08876* (2020) (cit. on pp. 5, 6).
- 426 [Bac14] Miroslav Bacák. “Convex analysis and optimization in Hadamard spaces”. In: *Convex*
427 *Analysis and Optimization in Hadamard Spaces*. de Gruyter, 2014 (cit. on pp. 2, 3, 6).
- 428 [BCO16] Glaydston de Carvalho Bento, João Xavier da Cruz-Neto, and Paulo Roberto Oliveira.
429 “A New Approach to the Proximal Point Method: Convergence on General Riemannian
430 Manifolds”. In: *J. Optim. Theory Appl.* 168.3 (2016), pp. 743–755 (cit. on p. 7).
- 431 [BFM17] Glaydston de Carvalho Bento, Orizon P. Ferreira, and Jefferson G. Melo. “Iteration-
432 Complexity of Gradient, Subgradient and Proximal Point Methods on Riemannian
433 Manifolds”. In: *J. Optim. Theory Appl.* 173.2 (2017), pp. 548–562 (cit. on pp. 1, 7).
- 434 [BFO15] GC Bento, OP Ferreira, and PR Oliveira. “Proximal point method for a special class of
435 nonconvex functions on Hadamard manifolds”. In: *Optimization* 64.2 (2015), pp. 289–
436 319 (cit. on p. 1).
- 437 [BHP13] A Barani, S Hosseini, and MR Pouryayevali. “On the metric projection onto φ -convex
438 subsets of Hadamard manifolds”. In: *Revista matemática complutense* 26.2 (2013),
439 pp. 815–826 (cit. on p. 6).

- 440 [BT09] Amir Beck and Marc Teboulle. “A Fast Iterative Shrinkage-Thresholding Algorithm for
441 [Linear Inverse Problems](#)”. In: *SIAM J. Imaging Sci.* 2.1 (2009), pp. 183–202 (cit. on
442 p. 2).
- 443 [CA16] Léopold Cambier and Pierre-Antoine Absil. “Robust Low-Rank Matrix Completion by
444 [Riemannian Optimization](#)”. In: *SIAM J. Scientific Computing* 38.5 (2016) (cit. on p. 1).
- 445 [Car+17] Yair Carmon, John C. Duchi, Oliver Hinder, and Aaron Sidford. ““Convex Until Proven
446 [Guilty](#)”: Dimension-Free Acceleration of Gradient Descent on Non-Convex Functions”.
447 In: *Proceedings of the 34th International Conference on Machine Learning, ICML 2017,*
448 *Sydney, NSW, Australia, 6-11 August 2017.* 2017, pp. 654–663 (cit. on p. 2).
- 449 [CB19] Chris Criscitiello and Nicolas Boumal. “Efficiently escaping saddle points on manifolds”.
450 In: *Advances in Neural Information Processing Systems 32: Annual Conference on*
451 *Neural Information Processing Systems 2019, NeurIPS 2019, 8-14 December 2019,*
452 *Vancouver, BC, Canada.* 2019, pp. 5985–5995 (cit. on p. 1).
- 453 [CB20] Chris Criscitiello and Nicolas Boumal. “An accelerated first-order method for non-
454 [convex optimization on manifolds](#)”. In: *arXiv preprint arXiv:2008.02252* (2020) (cit. on
455 p. 1).
- 456 [CB21] Christopher Criscitiello and Nicolas Boumal. “Negative curvature obstructs acceleration
457 [for geodesically convex optimization, even with exact first-order oracles](#)”. In: *CoRR*
458 *abs/2111.13263* (2021) (cit. on pp. 2, 5, 6).
- 459 [CDO18] Michael Cohen, Jelena Diakonikolas, and Lorenzo Orecchia. “On Acceleration with
460 [Noise-Corrupted Gradients](#)”. In: *Proceedings of the 35th International Conference on*
461 *Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15,*
462 *2018.* 2018, pp. 1018–1027 (cit. on p. 2).
- 463 [Cha+21] Shih-Sen Chang, Jen-Chih Yao, M Liu, and LC Zhao. “Inertial proximal point algorithm
464 [for variational inclusion in Hadamard manifolds](#)”. In: *Applicable Analysis* (2021), pp. 1–
465 12 (cit. on p. 7).
- 466 [CM12] Glaydston de Carvalho Bento and Jefferson G. Melo. “Subgradient Method for Convex
467 [Feasibility on Riemannian Manifolds](#)”. In: *J. Optim. Theory Appl.* 152.3 (2012), pp. 773–
468 785 (cit. on p. 1).
- 469 [CMP21] Francisco Criado, David Martínez-Rubio, and Sebastian Pokutta. “Fast Algorithms for
470 [Packing Proportional Fairness and its Dual](#)”. In: *arXiv preprint arXiv:2109.03678* (2021)
471 (cit. on p. 2).
- 472 [Cru+06] João Xavier da Cruz Neto, Orizon Pereira Ferreira, L. R. Lucambio Pérez, and Sándor
473 Zoltán Németh. “Convex- and Monotone-Transformable Mathematical Programming
474 [Problems and a Proximal-Like Point Method](#)”. In: *J. Glob. Optim.* 35.1 (2006), pp. 53–
475 69 (cit. on p. 1).
- 476 [CS17] Anoop Cherian and Suvrit Sra. “Riemannian Dictionary Learning and Sparse Coding
477 [for Positive Definite Matrices](#)”. In: *IEEE Trans. Neural Networks Learn. Syst.* 28.12
478 (2017), pp. 2859–2871 (cit. on p. 1).
- 479 [CS19] Ashok Cutkosky and Tamás Sarlós. “Matrix-Free Preconditioning in Online Learning”.
480 In: *Proceedings of the 36th International Conference on Machine Learning, ICML 2019,*
481 *9-15 June 2019, Long Beach, California, USA.* 2019, pp. 1455–1464 (cit. on p. 2).
- 482 [DJ19] Jelena Diakonikolas and Michael I. Jordan. “Generalized Momentum-Based Methods:
483 [A Hamiltonian Perspective](#)”. In: *CoRR abs/1906.00436* (2019) (cit. on p. 2).
- 484 [DN20] Nikita Doikov and Yurii E. Nesterov. “Contracting Proximal Methods for Smooth
485 [Convex Optimization](#)”. In: *SIAM J. Optim.* 30.4 (2020), pp. 3146–3169 (cit. on p. 2).
- 486 [DO18] Jelena Diakonikolas and Lorenzo Orecchia. “Accelerated Extra-Gradient Descent: A
487 [Novel Accelerated First-Order Method](#)”. In: *9th Innovations in Theoretical Computer*
488 *Science Conference, ITCS 2018, January 11-14, 2018, Cambridge, MA, USA.* 2018,
489 23:1–23:19 (cit. on p. 2).
- 490 [DO19] Jelena Diakonikolas and Lorenzo Orecchia. “The Approximate Duality Gap Technique:
491 [A Unified Theory of First-Order Methods](#)”. In: *SIAM Journal on Optimization* 29.1
492 (2019), pp. 660–689 (cit. on p. 2).
- 493 [DT14] Yoel Drori and Marc Teboulle. “Performance of first-order methods for smooth convex
494 [minimization: a novel approach](#)”. In: *Math. Program.* 145.1-2 (2014), pp. 451–482
495 (cit. on p. 2).

- 496 [EAS98] Alan Edelman, Tomás A. Arias, and Steven Thomas Smith. “The Geometry of Algo-
497 rithms with Orthogonality Constraints”. In: *SIAM J. Matrix Analysis Applications* 20.2
498 (1998), pp. 303–353 (cit. on p. 1).
- 499 [Gas+19] Alexander Gasnikov, Pavel E. Dvurechensky, Eduard A. Gorbunov, Evgeniya A.
500 Vorontsova, Daniil Selikhanovych, César A. Uribe, Bo Jiang, Haoyue Wang, Shuzhong
501 Zhang, Sébastien Bubeck, Qijia Jiang, Yin Tat Lee, Yuanzhi Li, and Aaron Sidford.
502 “Near Optimal Methods for Minimizing Convex Functions with Lipschitz p -th Deriva-
503 tives”. In: *Conference on Learning Theory, COLT 2019, 25-28 June 2019, Phoenix, AZ,
504 USA*. 2019, pp. 1392–1393 (cit. on p. 2).
- 505 [GHT15] Matthieu Genicot, Wen Huang, and Nickolay T. Trendafilov. “Weakly Correlated Sparse
506 Components with Nearly Orthonormal Loadings”. In: *Geometric Science of Information
507 - Second International Conference, GSI 2015, Palaiseau, France, October 28-30, 2015,
508 Proceedings*. 2015, pp. 484–490 (cit. on p. 1).
- 509 [HM21] Linus Hamilton and Ankur Moitra. “A No-go Theorem for Acceleration in the Hyper-
510 bolic Plane”. In: *arXiv preprint arXiv:2101.05657* (2021) (cit. on p. 5).
- 511 [HP13] Seyedehsomyeh Hosseini and Mohamad Pouryayevali. “On the metric projection
512 onto prox-regular subsets of Riemannian manifolds”. In: *Proceedings of the American
513 Mathematical Society* 141.1 (2013), pp. 233–244 (cit. on p. 6).
- 514 [HS15] Reshad Hosseini and Suvrit Sra. “Matrix Manifold Optimization for Gaussian Mixtures”.
515 In: *Advances in Neural Information Processing Systems 28: Annual Conference on
516 Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec,
517 Canada*. 2015, pp. 910–918 (cit. on p. 1).
- 518 [HS17] Reshad Hosseini and Suvrit Sra. “An Alternative to EM for Gaussian Mixture Models:
519 Batch and Stochastic Riemannian Optimization”. In: *CoRR abs/1706.03267* (2017)
520 (cit. on p. 1).
- 521 [HS18] Gennadij Heidel and Volker Schulz. “A Riemannian trust-region method for low-rank
522 tensor completion”. In: *Numerical Lin. Alg. with Applic.* 25.6 (2018) (cit. on p. 1).
- 523 [HS20] Reshad Hosseini and Suvrit Sra. “Recent advances in stochastic Riemannian optimiza-
524 tion”. In: *Handbook of Variational Methods for Nonlinear Geometric Data* (2020),
525 pp. 527–554 (cit. on p. 2).
- 526 [HSS19] Oliver Hinder, Aaron Sidford, and Nimit Sharad Sohoni. “Near-Optimal Methods for
527 Minimizing Star-Convex Functions and Beyond”. In: *CoRR abs/1906.11985* (2019)
528 (cit. on p. 2).
- 529 [HW19a] Wen Huang and Ke Wei. “Extending FISTA to Riemannian Optimization for Sparse
530 PCA”. In: *arXiv preprint arXiv:1909.05485* (2019) (cit. on p. 5).
- 531 [HW19b] Wen Huang and Ke Wei. “Riemannian Proximal Gradient Methods”. In: *arXiv preprint
532 arXiv:1909.06065* (2019) (cit. on p. 1).
- 533 [Iva+21] Anastasiya Ivanova, Dmitry Pasechnyuk, Dmitry Grishchenko, Egor Shulgin, Alexander
534 V. Gasnikov, and Vladislav Matyukhin. “Adaptive Catalyst for Smooth Convex Opti-
535 mization”. In: *Optimization and Applications - 12th International Conference, OPTIMA
536 2021, Petrovac, Montenegro, September 27 - October 1, 2021, Proceedings*. Ed. by
537 Nicholas N. Olenev, Yuri G. Evtushenko, Milojica Jacimovic, Michael Yu. Khachay,
538 and Vlasta Malkova. Vol. 13078. Lecture Notes in Computer Science. Springer, 2021,
539 pp. 20–37 (cit. on p. 2).
- 540 [Jou+20] Pooria Joulani, Anant Raj, András György, and Csaba Szepesvári. “A simpler approach
541 to accelerated optimization: iterative averaging meets optimism”. In: *Proceedings of
542 the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020,
543 Virtual Event*. Vol. 119. Proceedings of Machine Learning Research. PMLR, 2020,
544 pp. 4984–4993 (cit. on p. 2).
- 545 [JS21] Jikai Jin and Suvrit Sra. “A Riemannian Accelerated Proximal Extragradient Framework
546 and its Implications”. In: *CoRR abs/2111.02763* (2021) (cit. on p. 5).
- 547 [JTU03] Ian T Jolliffe, Nickolay T Trendafilov, and Mudassir Uddin. “A modified principal
548 component technique based on the LASSO”. In: *Journal of computational and Graphical
549 Statistics* 12.3 (2003), pp. 531–547 (cit. on p. 1).
- 550 [KG20] Dmitry Kamzolov and Alexander Gasnikov. “Near-optimal hyperfast second-order
551 method for convex optimization and its sliding”. In: *arXiv preprint arXiv:2002.09050*
552 (2020) (cit. on p. 2).

- 553 [Kha+21] Konrawut Khammahawong, Poom Kumam, Parin Chaipunya, and Juan Martínez-
554 Moreno. “Tseng’s methods for inclusion problems on Hadamard manifolds”. In: *Opti-*
555 *mization* (2021), pp. 1–35 (cit. on p. 7).
- 556 [KJM19] Hiroyuki Kasai, Pratik Jawanpuria, and Bamdev Mishra. “Riemannian adaptive stochas-
557 tic gradient algorithms on matrix manifolds”. In: *Proceedings of the 36th International*
558 *Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California,*
559 *USA*. 2019, pp. 3262–3271 (cit. on p. 1).
- 560 [KL17] Masoud Badiei Khuzani and Na Li. “Stochastic Primal-Dual Method on Riemannian
561 Manifolds of Bounded Sectional Curvature”. In: *16th IEEE International Conference on*
562 *Machine Learning and Applications, ICMLA 2017, Cancun, Mexico, December 18-21,*
563 *2017*. 2017, pp. 133–140 (cit. on p. 2).
- 564 [KY22] Jungbin Kim and Insoon Yang. “Nesterov Acceleration for Riemannian Optimization”.
565 In: *arXiv preprint arXiv:2202.02036* (2022) (cit. on pp. 1, 2, 5, 6, 8, 22).
- 566 [Lev77] Tullio Levi-Civita. *The absolute differential calculus (calculus of tensors)*. Courier
567 Corporation, 1977. ISBN: 978-0-486-31625-3 (cit. on p. 3).
- 568 [Lez20] Mario Lezcano-Casado. “Curvature-Dependant Global Convergence Rates for Optimiza-
569 tion on Manifolds of Bounded Geometry”. In: *arXiv preprint arXiv:2008.02517* (2020)
570 (cit. on p. 4).
- 571 [Lin+20] Lizhen Lin, Bayan Saparbayeva, Michael Minyi Zhang, and David B. Dunson. “Accel-
572 erated Algorithms for Convex and Non-Convex Optimization on Manifolds”. In: *CoRR*
573 *abs/2010.08908* (2020) (cit. on p. 5).
- 574 [Liu+17] Yuanyuan Liu, Fanhua Shang, James Cheng, Hong Cheng, and Licheng Jiao. “Ac-
575 celerated First-order Methods for Geodesically Convex Optimization on Riemannian
576 Manifolds”. In: *Advances in Neural Information Processing Systems 30: Annual Confer-*
577 *ence on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach,*
578 *CA, USA*. Ed. by Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach,
579 Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett. 2017, pp. 4868–4877 (cit. on
580 p. 5).
- 581 [LM19] Mario Lezcano-Casado and David Martínez-Rubio. “Cheap Orthogonal Constraints in
582 Neural Networks: A Simple Parametrization of the Orthogonal and Unitary Group”. In:
583 *Proceedings of the 36th International Conference on Machine Learning, ICML 2019,*
584 *9-15 June 2019, Long Beach, California, USA*. 2019, pp. 3794–3803 (cit. on p. 1).
- 585 [LMH17] Hongzhou Lin, Julien Mairal, and Zaid Harchaoui. “Catalyst Acceleration for First-order
586 Convex Optimization: from Theory to Practice”. In: *J. Mach. Learn. Res.* 18 (2017),
587 212:1–212:54 (cit. on p. 2).
- 588 [Mar21] David Martínez-Rubio. “Acceleration in first-order optimization methods: promenading
589 beyond convexity or smoothness, and applications”. PhD thesis. University of Oxford,
590 2021 (cit. on p. 2).
- 591 [Mar22] David Martínez-Rubio. “Global Riemannian Acceleration in Hyperbolic and Spherical
592 Spaces”. In: *International Conference on Algorithmic Learning Theory, 29-1 April 2022,*
593 *Paris, France*. Ed. by Sanjoy Dasgupta and Nika Haghtalab. Vol. 167. Proceedings of
594 Machine Learning Research. PMLR, 2022, pp. 768–826 (cit. on pp. 4–6, 9, 21).
- 595 [MS14] Bamdev Mishra and Rodolphe Sepulchre. “R3MC: A Riemannian three-factor algorithm
596 for low-rank matrix completion”. In: *53rd IEEE Conference on Decision and Control,*
597 *CDC 2014, Los Angeles, CA, USA, December 15-17, 2014*. 2014, pp. 1137–1142 (cit. on
598 p. 1).
- 599 [Nes05] Yurii Nesterov. “Smooth minimization of non-smooth functions”. In: *Math. Program.*
600 103.1 (2005), pp. 127–152 (cit. on pp. 4, 6).
- 601 [Nes83] Yurii Nesterov. “A method of solving a convex programming problem with convergence
602 rate $O(1/k^2)$ ”. In: *Soviet Mathematics Doklady*. Vol. 27. 1983, pp. 372–376 (cit. on pp. 2,
603 4).
- 604 [Pet06] Peter Petersen. *Riemannian geometry*. Ed. by S Axler and KA Ribet. Vol. 171. Springer,
605 2006. ISBN: 978-0-387-29403-2 (cit. on p. 2).
- 606 [SBC16] Weijie Su, Stephen P. Boyd, and Emmanuel J. Candès. “A Differential Equation for
607 Modeling Nesterov’s Accelerated Gradient Method: Theory and Insights”. In: *J. Mach.*
608 *Learn. Res.* 17 (2016), 153:1–153:43 (cit. on p. 2).

- 609 [SFF19] Yue Sun, Nicolas Flammarion, and Maryam Fazel. “Escaping from saddle points on
610 Riemannian manifolds”. In: *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, 8-14 December 2019, Vancouver, BC, Canada*. 2019, pp. 7274–7284 (cit. on p. 1).
- 613 [SKM17] Hiroyuki Sato, Hiroyuki Kasai, and Bamdev Mishra. “Riemannian stochastic variance
614 reduced gradient”. In: *CoRR* abs/1702.05594 (2017) (cit. on p. 2).
- 615 [SKM19] Hiroyuki Sato, Hiroyuki Kasai, and Bamdev Mishra. “Riemannian Stochastic Variance
616 Reduced Gradient Algorithm with Retraction and Vector Transport”. In: *SIAM Journal on Optimization* 29.2 (2019), pp. 1444–1472 (cit. on p. 2).
- 618 [SQW17] Ju Sun, Qing Qu, and John Wright. “Complete Dictionary Recovery Over the Sphere
619 II: Recovery by Riemannian Trust-Region Method”. In: *IEEE Trans. Inf. Theory* 63.2
620 (2017), pp. 885–914 (cit. on p. 1).
- 621 [Tan+14] Mingkui Tan, Ivor W. Tsang, Li Wang, Bart Vandereycken, and Sinno Jialin Pan. “Rie-
622 mannian Pursuit for Big Matrix Recovery”. In: *Proceedings of the 31th International
623 Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*. 2014,
624 pp. 1539–1547 (cit. on p. 1).
- 625 [TH14] Guo-ji Tang and Nan-Jing Huang. “Rate of convergence for proximal point algorithms
626 on Hadamard manifolds”. In: *Oper. Res. Lett.* 42.6-7 (2014), pp. 383–387 (cit. on p. 7).
- 627 [Tri+18] Nilesh Tripurani, Nicolas Flammarion, Francis Bach, and Michael I. Jordan. “Averag-
628 ing Stochastic Gradient Descent on Riemannian Manifolds”. In: *CoRR* abs/1802.09128
629 (2018) (cit. on p. 2).
- 630 [Tse08] Paul Tseng. “On accelerated proximal gradient methods for convex-concave optimiza-
631 tion”. In: *submitted to SIAM Journal on Optimization* 2.3 (2008) (cit. on p. 2).
- 632 [Van13] Bart Vandereycken. “Low-Rank Matrix Completion by Riemannian Optimization”. In:
633 *SIAM Journal on Optimization* 23.2 (2013), pp. 1214–1236 (cit. on p. 1).
- 634 [Wal74] Rolf Walter. “On the metric projection onto convex sets in Riemannian spaces”. In:
635 *Archiv der Mathematik* 25.1 (1974), pp. 91–98 (cit. on p. 6).
- 636 [Wan+15] Jinhua Wang, Chong Li, Genaro López-Acedo, and Jen-Chih Yao. “Convergence analy-
637 sis of inexact proximal point algorithms on Hadamard manifolds”. In: *J. Glob. Optim.*
638 61.3 (2015), pp. 553–573 (cit. on p. 7).
- 639 [Wan+16] Jinhua Wang, Chong Li, Genaro López-Acedo, and Jen-Chih Yao. “Proximal Point
640 Algorithms on Hadamard Manifolds: Linear Convergence and Finite Termination”. In:
641 *SIAM J. Optim.* 26.4 (2016), pp. 2696–2729 (cit. on p. 7).
- 642 [Wei+16] Ke Wei, Jian-Feng Cai, Tony F Chan, and Shingyu Leung. “Guarantees of Riemannian
643 optimization for low rank matrix completion”. In: *arXiv preprint arXiv:1603.06610*
644 (2016) (cit. on p. 1).
- 645 [Wie12] Ami Wiesel. “Geodesic Convexity and Covariance Estimation”. In: *IEEE Trans. Signal
646 Process.* 60.12 (2012), pp. 6182–6189 (cit. on p. 1).
- 647 [WRM16] Di Wang, Satish Rao, and Michael W. Mahoney. “Unified Acceleration Method for
648 Packing and Covering Problems via Diameter Reduction”. In: *43rd International Col-
649 loquium on Automata, Languages, and Programming, ICALP 2016, July 11-15, 2016,
650 Rome, Italy*. 2016, 50:1–50:13 (cit. on p. 2).
- 651 [WS17] Melanie Weber and Suvrit Sra. “Frank-Wolfe methods for geodesically convex optimiza-
652 tion with application to the matrix geometric mean”. In: *CoRR* abs/1710.10770 (2017)
653 (cit. on p. 1).
- 654 [WS19] Melanie Weber and Suvrit Sra. “Nonconvex stochastic optimization on manifolds via
655 Riemannian Frank-Wolfe methods”. In: *CoRR* abs/1910.04194 (2019) (cit. on p. 1).
- 656 [WWJ16] Andre Wibisono, Ashia C. Wilson, and Michael I. Jordan. “A Variational Perspective on
657 Accelerated Methods in Optimization”. In: *CoRR* abs/1603.04245 (2016) (cit. on p. 2).
- 658 [ZRS16] Hongyi Zhang, Sashank J. Reddi, and Suvrit Sra. “Riemannian SVRG: Fast Stochastic
659 Optimization on Riemannian Manifolds”. In: *Advances in Neural Information Process-
660 ing Systems 29: Annual Conference on Neural Information Processing Systems 2016,
661 December 5-10, 2016, Barcelona, Spain*. 2016, pp. 4592–4600 (cit. on p. 2).
- 662 [ZS16] Hongyi Zhang and Suvrit Sra. “First-order Methods for Geodesically Convex Optimiza-
663 tion”. In: *Proceedings of the 29th Conference on Learning Theory, COLT 2016, New
664 York, USA, June 23-26, 2016*. 2016, pp. 1617–1638 (cit. on pp. 1, 2, 6, 9, 21, 24, 25).

- 665 [ZS18] Hongyi Zhang and Suvrit Sra. “An Estimate Sequence for Geodesically Convex Opti-
666 mization”. In: *Conference On Learning Theory, COLT 2018, Stockholm, Sweden, 6-9*
667 *July 2018*. Ed. by Sébastien Bubeck, Vianney Perchet, and Philippe Rigollet. Vol. 75.
668 *Proceedings of Machine Learning Research*. PMLR, 2018, pp. 1703–1723 (cit. on pp. 5,
669 6, 8).
- 670 [ZYF19] Pan Zhou, Xiao-Tong Yuan, and Jiashi Feng. “Faster First-Order Methods for Stochastic
671 Non-Convex Optimization on Riemannian Manifolds”. In: *The 22nd International*
672 *Conference on Artificial Intelligence and Statistics, AISTATS 2019, 16-18 April 2019,*
673 *Naha, Okinawa, Japan*. 2019, pp. 138–147 (cit. on p. 1).
- 674 [ZZS18] Jingzhao Zhang, Hongyi Zhang, and Suvrit Sra. “R-SPIDER: A Fast Riemannian
675 Stochastic Optimization Algorithm with Curvature Independent Rate”. In: *CoRR*
676 *abs/1811.04194* (2018) (cit. on p. 1).
- 677 [ZZS22] Peiyuan Zhang, Jingzhao Zhang, and Suvrit Sra. “Minimax in Geodesic Metric Spaces:
678 *Sion’s Theorem and Algorithms*”. In: *CoRR abs/2202.06950* (2022) (cit. on p. 2).

679 Checklist

- 680 1. For all authors...
- 681 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
682 contributions and scope? [Yes]
- 683 (b) Did you describe the limitations of your work? [Yes] See [Section 3](#).
- 684 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 685 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
686 them? [Yes]
- 687 2. If you are including theoretical results...
- 688 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See the begin-
689 ning of [Section 1.3](#). Alternatively, see the statements of [Theorem 2.2](#) and [Theorem 2.4](#).
- 690 (b) Did you include complete proofs of all theoretical results? [Yes]
- 691 3. If you ran experiments...
- 692 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
693 mental results (either in the supplemental material or as a URL)? [N/A]
- 694 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
695 were chosen)? [N/A]
- 696 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
697 ments multiple times)? [N/A]
- 698 (d) Did you include the total amount of compute and the type of resources used (e.g., type
699 of GPUs, internal cluster, or cloud provider)? [N/A]
- 700 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 701 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 702 (b) Did you mention the license of the assets? [N/A]
- 703 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 704
- 705 (d) Did you discuss whether and how consent was obtained from people whose data you’re
706 using/curating? [N/A]
- 707 (e) Did you discuss whether the data you are using/curating contains personally identifiable
708 information or offensive content? [N/A]
- 709 5. If you used crowdsourcing or conducted research with human subjects...
- 710 (a) Did you include the full text of instructions given to participants and screenshots, if
711 applicable? [N/A]
- 712 (b) Did you describe any potential participant risks, with links to Institutional Review
713 Board (IRB) approvals, if applicable? [N/A]
- 714 (c) Did you include the estimated hourly wage paid to participants and the total amount
715 spent on participant compensation? [N/A]

716 **A Optimization lemmas and proofs**

717 We start by noting a property that our parameters satisfy.

718 **Lemma A.1.** *For the parameter choices of a_k and A_{k-1} in Algorithm 1 we have, for all $k \geq 1$:*

$$\frac{8\lambda}{9}(\xi A_{k-1} + a_k) \geq a_k^2 \geq \frac{3\lambda}{4}(\xi A_{k-1} + \xi a_k).$$

719 *Proof.* It is a simple computation to check that a_k and A_{k-1} satisfy such inequality. The inequalities
720 are equivalent to the following, which trivially holds:

$$\begin{aligned} \frac{8}{9}((k^2 - k + 64k\xi - 64\xi + 1000\xi^2) + (2k + 64\xi)) &\geq \frac{4}{5}(k^2 + 64k\xi + 1024\xi^2) \\ &\geq \frac{3}{4}((k^2 - k + 64k\xi - 64\xi + 1000\xi^2) + (2k\xi + 64\xi^2)) \end{aligned}$$

721

□

722 We now prove Proposition 2.1, which will allow us to use ψ_k as a Lyapunov function to show the
723 final convergence rates. The proof will use Lemma A.2 and Lemma A.3, that we state and prove
724 afterwards.

725 *Proof (Proposition 2.1).* Inequality $(1 - \Delta_k)\psi_k \leq \psi_{k-1}$ is equivalent to

$$\begin{aligned} (1 - \Delta_k) \left(A_k(f(y_k) - f(x^*)) + \frac{1}{2}\|z_k^{y_k} - x^*\|_{y_k}^2 + \frac{\xi - 1}{2}\|y_k - z_k^{y_k}\|_{y_k}^2 \right) \\ \leq A_{k-1}(f(y_{k-1}) - f(x^*)) + \left(\frac{1}{2}\|z_{k-1}^{y_{k-1}} - x^*\|_{y_{k-1}}^2 + \frac{\xi - 1}{2}\|y_{k-1} - z_{k-1}^{y_{k-1}}\|_{y_{k-1}}^2 \right) \end{aligned}$$

726 which, by bounding $(1 - \Delta_k)(f(y_k) - f(x^*)) \leq f(y_k) - f(x^*)$ and reorganizing, is implied by the
727 following:

$$\begin{aligned} A_{k-1}(f(y_k) - f(y_{k-1})) + \frac{a_k}{\xi}(f(y_k) - f(x^*)) &\leq \frac{1}{2}\|z_{k-1}^{y_{k-1}} - x^*\|_{y_{k-1}}^2 - \frac{1 - \Delta_k}{2}\|z_k^{y_k} - x^*\|_{y_k}^2 \\ &+ \frac{\xi - 1}{2} \left(\|y_{k-1} - z_{k-1}^{y_{k-1}}\|_{y_{k-1}}^2 - (1 - \Delta_k)\|y_k - z_k^{y_k}\|_{y_k}^2 \right). \end{aligned}$$

728 We have that due to the projection in Line 12, then x_k is not very far from any $p \in \mathcal{X}$:

$$d(x_k, p) \leq \|x_k - y_{k-1}\|_{y_{k-1}} + d(y_{k-1}, p) \stackrel{\textcircled{1}}{<} \|\bar{z}_{k-1}^{y_{k-1}} - y_{k-1}\|_{y_{k-1}} + D \stackrel{\textcircled{2}}{\leq} 2D, \quad (1)$$

729 where $\textcircled{1}$ holds by the definition of x_k and the fact $y_{k-1}, p \in \mathcal{X}$, and $\textcircled{2}$ is due to the projection
730 defining $\bar{z}_{k-1}^{y_{k-1}}$. Now we use the first part of Lemma A.2 with both $x \leftarrow y_{k-1}$ and $x \leftarrow x^*$ and we
731 bound the resulting errors $\varepsilon_k(\cdot)$ by using the second part of Lemma A.2. We also use Lemma A.3, so
732 it is enough to prove the following:

$$\begin{aligned} A_{k-1}\langle v_k^x, x_k - y_{k-1} \rangle + (a_k/\xi)\langle v_k^x, x_k - z_{k-1}^{x_k} + z_{k-1}^{x_k} - x^* \rangle - \frac{4\lambda}{9}(A_{k-1} + a_k/\xi)\|v_k^x\|^2 \\ \leq \frac{1}{2}\|z_{k-1}^{x_k} - x^*\|_{x_k}^2 - \frac{1}{2}\|z_k^{x_k} - x^*\|_{x_k}^2 + \frac{\xi - 1}{2} \left(\|x_k - z_{k-1}^{x_k}\|_{x_k}^2 - \|x_k - z_k^{x_k}\|_{x_k}^2 \right), \end{aligned}$$

733 Note that thanks to Lemma A.3 now we have the potentials on the right hand side as expressions in
734 the tangent space of x_k . Also, note that we have canceled some potentials proportional to Δ_k coming
735 from the bound on the error $\varepsilon_k(\cdot)$. Now we use that by definition of x_k we have, for all $v \in T_{x_k} \mathcal{M}$,
736 $A_{k-1}\langle v, x_k - y_{k-1} \rangle = -a_k\langle v, x_k - z_{k-1}^{x_k} \rangle$, so we use this fact for $v = v_k^x$ and use the following
737 identity, that holds by the definition of $z_k^{x_k} \stackrel{\text{def}}{=} z_{k-1}^{x_k} - \eta_k v_k^x$:

$$\frac{a_k/\xi}{\eta_k}\langle \eta_k v_k^x, z_{k-1}^{x_k} - x^* \rangle = \frac{a_k/\xi}{2\eta_k} \left(\eta_k^2 \|v_k^x\|_{x_k}^2 + \|z_{k-1}^{x_k} - x^*\|_{x_k}^2 - \|z_k^{x_k} - x^*\|_{x_k}^2 \right).$$

738 so that, after canceling terms, it is enough to prove:

$$\begin{aligned}
& a_k(1 - 1/\xi)\langle -v_k^x, x_k - z_{k-1}^{x_k} \rangle - \frac{a_k(1 - 1/\xi)}{2\eta_k} \eta_k^2 \|v_k^x\|^2 \\
& + \|v_k^x\|^2 \left(-\frac{4}{9}(A_{k-1}\lambda + a_k\lambda/\xi) + \frac{a_k\eta_k}{2} \right) \\
& \leq \frac{\xi - 1}{2} \left(\|x_k - z_{k-1}^{x_k}\|_{x_k}^2 - \|x_k - z_k^{x_k}\|_{x_k}^2 \right),
\end{aligned} \tag{2}$$

739 Now we show that in the previous inequality (2), the first line cancels with the last line. Note that
740 $(a_k(1 - 1/\xi))/\eta_k = (1 - 1/\xi)/(1/\xi) = \xi - 1$. Thus, by using again the definition of $z_k^{x_k}$, we have:

$$\frac{a_k(1 - 1/\xi)}{\eta_k} \langle -\eta_k v_k^x, x_k - z_{k-1}^{x_k} \rangle = \frac{a_k(1 - 1/\xi)}{2\eta_k} \left(\eta_k^2 \|v_k^x\|_{x_k}^2 + \|x_k - z_{k-1}^{x_k}\|_{x_k}^2 - \|x_k - z_k^{x_k}\|_{x_k}^2 \right).$$

741 Finally, it only remains to prove:

$$\frac{\|v_k^x\|^2}{2\xi} \cdot \left(-\frac{8}{9}(\xi A_{k-1}\lambda + a_k\lambda) + a_k^2 \right) \leq 0, \tag{3}$$

742 which holds by [Lemma A.1](#). □

743 We now show the two auxiliary lemmas that we used in the previous proof.

744 **Lemma A.2.** Let $h_k(x) \stackrel{\text{def}}{=} f(x) + \frac{1}{2\lambda}d(x_k, x)^2$ be the strongly g -convex function used at step k , and
745 let $y_k^* = \arg \min_{y \in \mathcal{X}} h_k(y)$. Then, for $y_k \in \mathcal{X}$, if we let $v_k^x \stackrel{\text{def}}{=} -\text{Log}_{x_k}(y_k)/\lambda$, then the following
746 holds, for all $x \in \mathcal{X}$:

$$f(x) \geq f(y_k) + \langle v_k^x, x - x_k \rangle_{x_k} + \frac{\lambda}{2} \|v_k^x\|^2 - \varepsilon_k(x)$$

747 where $\varepsilon_k(x) \stackrel{\text{def}}{=} -\frac{1}{\lambda} \langle y_k - y_k^*, x - x_k \rangle_{x_k} + (h_k(y_k) - h_k(y_k^*))$. Moreover, if y_k satisfies

$$h_k(y_k) - h_k(y_k^*) \leq \frac{\Delta_k d(x_k, y_k^*)^2}{78\lambda},$$

748 then we have

$$\begin{aligned}
& -\frac{\lambda}{2} \|v_k^x\|^2 (A_{k-1} + a_k/\xi) + a_k \varepsilon_k(x^*)/\xi + A_{k-1} \varepsilon_k(y_{k-1}) \\
& \leq -\frac{4\lambda \|v_k^x\|^2}{9} (A_{k-1} + a_k/\xi) + \frac{\Delta_k}{2} \left(\|x^* - z_{k-1}^{x_k}\|_{x_k}^2 + (\xi - 1) \|x_k - z_{k-1}^{x_k}\|_{x_k}^2 \right).
\end{aligned}$$

749 *Proof.* The function h_k is $\frac{1}{\lambda}$ -strongly g -convex because by [Fact 1.3](#) the function $\frac{1}{2}d(x_k, x)^2$ is 1-
750 strongly g -convex in a Hadamard manifold. By the first-order optimality condition of h_k at y_k^* we
751 have that $\tilde{v}_k^y \stackrel{\text{def}}{=} \lambda^{-1} \text{Log}_{y_k^*}(x_k) \in \partial(f + I_{\mathcal{X}})(y_k^*)$ is a subgradient of $f + I_{\mathcal{X}}$ at y_k^* . Thus, we have,
752 for all $x \in \mathcal{X}$ and for $\tilde{v}_k^x \stackrel{\text{def}}{=} \Gamma_{y_k^*}^{x_k}(\tilde{v}_k^y)$:

$$\begin{aligned}
f(x) & \stackrel{\textcircled{1}}{\geq} f(y_k^*) + \langle \tilde{v}_k^y, x - y_k^* \rangle_{y_k^*} \\
& \stackrel{\textcircled{2}}{\geq} f(y_k^*) + \langle \tilde{v}_k^x, x - x_k \rangle_{x_k} + \lambda \|\tilde{v}_k^x\|^2 \\
& \stackrel{\textcircled{3}}{=} f(y_k) + \langle v_k^x, x - x_k \rangle_{x_k} + \frac{\lambda}{2} \|v_k^x\|^2 + \frac{\lambda}{2} \|\tilde{v}_k^x\|^2 \\
& + \langle \tilde{v}_k^x - v_k^x, x - x_k \rangle_{x_k} + \left((f(y_k^*) + \frac{\lambda}{2} \|\tilde{v}_k^x\|^2) - (f(y_k) + \frac{\lambda}{2} \|v_k^x\|^2) \right) \\
& \stackrel{\textcircled{4}}{\geq} f(y_k) + \langle v_k^x, x - x_k \rangle_{x_k} + \frac{\lambda}{2} \|v_k^x\|^2 + \frac{1}{\lambda} \langle y_k - y_k^*, x - x_k \rangle_{x_k} - (h_k(y_k) - h_k(y_k^*))
\end{aligned}$$

753 where ① holds because $\tilde{v}_k^y \in \partial(f + I_{\mathcal{X}})(y_k^*)$ and $x, y_k^* \in \mathcal{X}$. In ②, we used the first part of
 754 Lemma B.5 along with $\delta = 1$. We just added and subtracted some terms in ③, and in ④, we dropped
 755 $\frac{\lambda}{2} \|\tilde{v}_k^x\|^2$, and we used the definitions of h_k , \tilde{v}_k^x , and $v_k^x = -\text{Log}_{x_k}(y_k)/\lambda$.

756 Now we proceed to prove the second part. The following holds:

$$\begin{aligned}
 & -\frac{a_k}{\lambda\xi} \langle y_k - y_k^*, x^* - x_k \rangle_{x_k} - A_{k-1} \frac{1}{\lambda} \langle y_k - y_k^*, y_{k-1} - x_k \rangle_{x_k} \\
 & \stackrel{\textcircled{1}}{\leq} \frac{1}{\lambda} \|y_k - y_k^*\|_{x_k} \cdot \left\| \frac{a_k}{\xi} x^* + A_{k-1} y_{k-1} - \left(\frac{a_k}{\xi} + A_{k-1} \right) x_k \right\|_{x_k} \\
 & \stackrel{\textcircled{2}}{\leq} \frac{1}{\lambda} d(y_k, y_k^*) \cdot \frac{a_k}{\xi} \|x^* - z_{k-1}^{x_k} + (\xi - 1)(x_k - z_{k-1}^{x_k})\|_{x_k} \\
 & \stackrel{\textcircled{3}}{\leq} \frac{1}{\lambda} \sqrt{2\lambda(h_k(y_k) - h_k(y_k^*))} \cdot \frac{a_k}{\xi} \sqrt{\xi} \sqrt{\|x^* - z_{k-1}^{x_k}\|_{x_k}^2 + (\xi - 1)\|(x_k - z_{k-1}^{x_k})\|_{x_k}^2} \\
 & = \sqrt{\frac{2a_k^2(h_k(y_k) - h_k(y_k^*))}{\Delta_k \lambda \xi}} \cdot \sqrt{\Delta_k} \sqrt{\|x^* - z_{k-1}^{x_k}\|_{x_k}^2 + (\xi - 1)\|(x_k - z_{k-1}^{x_k})\|_{x_k}^2} \\
 & \stackrel{\textcircled{4}}{\leq} \frac{a_k^2(h_k(y_k) - h_k(y_k^*))}{\Delta_k \lambda \xi} + \frac{\Delta_k}{2} (\|x^* - z_{k-1}^{x_k}\|_{x_k}^2 + (\xi - 1)\|(x_k - z_{k-1}^{x_k})\|_{x_k}^2),
 \end{aligned} \tag{4}$$

757 where ① groups some terms and uses Cauchy-Schwartz. In inequality ②, for the first term we
 758 bounded the distance between y_k^* and y_k estimated from $T_{x_k} \mathcal{M}$ by the actual distance, which is a
 759 property that holds in Hadamard manifolds and it holds by the first part of Corollary B.2 with $\delta = 1$,
 760 $p \leftarrow y_k^*$, $y \leftarrow y_k$, $x \leftarrow x_k$, $z^y \leftarrow 0$. The second term is substituted by a term of equal value by
 761 using Euclidean trigonometry in $T_{x_k} \mathcal{M}$, as in the following. Let $w \stackrel{\text{def}}{=} \frac{1}{a_k/\xi + A_{k-1}} \left(\frac{a_k}{\xi} \text{Log}_{x_k}(x^*) + \right.$
 762 $\left. A_{k-1} \text{Log}_{x_k}(y_{k-1}) \right)$ and let $u \in T_{x_k}$ be the point in the line containing $\text{Log}_{x_k}(y_{k-1})$ and $0 =$
 763 $\text{Log}_{x_k}(x_k) \in T_{x_k}$ such that the triangle with vertices 0 , $\text{Log}_{x_k}(y_{k-1})$ and w and the triangle with
 764 vertices u , $\text{Log}_{x_k}(y_{k-1})$ and $\text{Log}_{x_k}(x^*)$ are similar triangles, and so

$$\frac{\|\text{Log}_{x_k}(x^*) - u\|}{\|w - \text{Log}_{x_k}(x_k)\|} \stackrel{\textcircled{5}}{=} \frac{\|\text{Log}_{x_k}(x^*) - \text{Log}_{x_k}(y_{k-1})\|}{\|w - \text{Log}_{x_k}(y_{k-1})\|} \stackrel{\textcircled{6}}{=} \frac{A_{k-1} + a_k/\xi}{a_k/\xi}. \tag{5}$$

765 We used the triangle similarity in ⑤ and in ⑥ we used the definition of w as a convex combination
 766 of $\text{Log}_{x_k}(x^*)$ and $\text{Log}_{x_k}(y_{k-1})$. It is enough to show $u = \xi z_{k-1}^{x_k}$ as in such a case what we applied
 767 in ② is equivalent to the equality (5) above. By the definition of x_k , we have ⑦ below and by
 768 triangle similarity we have ⑧ below:

$$z_{k-1}^{x_k} \stackrel{\textcircled{7}}{=} -\frac{A_{k-1}}{a_k} \text{Log}_{x_k}(y_{k-1}) \stackrel{\textcircled{8}}{=} \frac{A_{k-1}}{a_k} \cdot \frac{a_k/\xi}{A_{k-1}} u = \frac{1}{\xi} u,$$

769 as desired. In the next inequality ③, we used that by $(1/\lambda)$ -strong g-convexity of h_k and by optimality
 770 of y_k^* , we have $\frac{1}{2\lambda} d(\cdot, y_k^*)^2 \leq h_k(\cdot) - h_k(y_k^*)$. For the second term, we used that for vectors $b, c \in \mathbb{R}^n$
 771 and $\omega \in \mathbb{R}_{\geq 0}$, we have, by Young's inequality, $\|b + \omega c\| = \sqrt{\|b\|^2 + \omega^2 \|c\|^2 + 2\langle \sqrt{\omega} b, \sqrt{\omega} c \rangle} \leq$
 772 $\sqrt{(1 + \omega)(\|b\|^2 + \omega \|c\|^2)}$. In ④ we used Young's inequality.

773 Before we conclude, we note that

$$d(x_k, y_k^*) \leq \sqrt{2} d(x_k, y_k), \tag{6}$$

774 which is implied by the following, where we use the same as in ③ above, the assumption on y_k and
 775 $\Delta_k \leq 1$:

$$\begin{aligned}
 d(x_k, y_k^*) & \leq d(x_k, y_k) + d(y_k, y_k^*) \leq d(x_k, y_k) + \sqrt{2\lambda(h_k(y_k) - h_k(y_k^*))} \\
 & \leq d(x_k, y_k) + \sqrt{\Delta_k/34} \cdot d(x_k, y_k^*) \leq d(x_k, y_k) + d(x_k, y_k^*)/4.
 \end{aligned}$$

776 Finally, we can make use of (4) and (6) to obtain the claim in the second part of the lemma:

$$\begin{aligned}
& -\frac{\lambda}{2}\|v_k^x\|^2(A_{k-1} + a_k/\xi) + a_k\varepsilon_k(x^*)/\xi + A_{k-1}\varepsilon_k(y_{k-1}) - \frac{\Delta_k}{2}\|x^* - z_{k-1}^{x_k}\|_{x_k}^2 \\
& \quad - \Delta_k \frac{\xi - 1}{2} \|(x_k - z_{k-1}^{x_k})\|_{x_k}^2 \\
& \leq -\frac{\lambda}{2}\|v_k^x\|^2(A_{k-1} + a_k/\xi) + \left(A_{k-1} + a_k/\xi + \frac{a_k^2}{\Delta_k \lambda \xi}\right) (h_k(y_k) - h_k(y_k^*)) \\
& \stackrel{\textcircled{1}}{\leq} -\frac{\lambda}{2}\|v_k^x\|^2(A_{k-1} + a_k/\xi) + (A_{k-1} + a_k/\xi) \left(1 + \frac{a_k^2}{(\xi A_{k-1} + a_k)\lambda}\right) \frac{d(x_k, y_k)^2}{34\lambda} \\
& \stackrel{\textcircled{2}}{\leq} -\frac{\lambda}{2}\|v_k^x\|^2(A_{k-1} + a_k/\xi) + \frac{d(x_k, y_k)^2}{18\lambda} (A_{k-1} + a_k/\xi) \\
& \stackrel{\textcircled{3}}{=} -\frac{4\lambda\|v_k^x\|^2}{9} (A_{k-1} + a_k/\xi),
\end{aligned}$$

777 where $\textcircled{1}$ holds by the assumption on y_k , $\Delta_k \leq 1$, and (6). Inequality $\textcircled{2}$ uses the upper bound on a_k^2
778 in Lemma A.1, and $\textcircled{3}$ uses the definition $v_k^x \stackrel{\text{def}}{=} -\text{Log}_{x_k}(y_k)/\lambda$.

779 □

780 The following lemma allows to *move* the regularized lower bounds on the objective without incurring
781 extra geometric penalties.

782 **Lemma A.3 (Translating Potentials with no Geometric Penalty).** *Using the variables in Algo-*
783 *rithm 1, for any $\Delta_k \in [0, 1]$, we have*

$$\begin{aligned}
& \|z_{k-1}^{x_k} - x^*\|_{x_k}^2 - (1 - \Delta_k)\|z_k^{x_k} - x^*\|_{x_k}^2 + (\xi - 1) \left(\|x_k - z_{k-1}^{x_k}\|_{x_k}^2 - (1 - \Delta_k)\|x_k - z_k^{x_k}\|_{x_k}^2 \right) \\
& \leq \|z_{k-1}^{y_{k-1}} - x^*\|_{y_{k-1}}^2 - (1 - \Delta_k)\|z_k^{y_k} - x^*\|_{y_k}^2 \\
& \quad + (\xi - 1) \left(\|y_{k-1} - z_{k-1}^{y_{k-1}}\|_{y_{k-1}}^2 - (1 - \Delta_k)\|y_k - z_k^{y_k}\|_{y_k}^2 \right).
\end{aligned}$$

784 *Proof.* Firstly, by the projection step in Line 12, we have

$$\|z_{k-1}^{y_{k-1}} - x^*\|_{y_k}^2 \geq \|\bar{z}_{k-1}^{y_{k-1}} - x^*\|_{y_k}^2 \quad \text{and} \quad (\xi - 1)\|z_{k-1}^{y_{k-1}}\|_{y_k}^2 \geq (\xi - 1)\|\bar{z}_{k-1}^{y_{k-1}}\|_{y_k}^2 \quad (7)$$

785 since the operation is a simple Euclidean projection onto the closed ball $\bar{B}(0, D)$ in $T_{y_k}\mathcal{M}$. By the
786 second part of Corollary B.2, $y = x_k$ and $x = y_{k-1}$ and by (1), we have $\textcircled{1}$ below

$$\begin{aligned}
& \|\bar{z}_{k-1}^{y_{k-1}} - x^*\|_{y_{k-1}}^2 + (\xi - 1)\|\bar{z}_{k-1}^{y_{k-1}}\|_{y_{k-1}}^2 \stackrel{\textcircled{1}}{\geq} \|z_{k-1}^{x_k} - x^*\|_{x_k}^2 + (\zeta_{2D} - 1)\|z_{k-1}^{x_k}\|_{x_k}^2 + (\xi - \zeta_{2D})\|z_{k-1}^{y_{k-1}}\|_{y_{k-1}}^2 \\
& \stackrel{\textcircled{2}}{\geq} \|z_{k-1}^{x_k} - x^*\|_{x_k}^2 + (\xi - 1)\|z_{k-1}^{x_k}\|_{x_k}^2 + (\xi - \zeta_{2D}) \left(\left(\frac{A_{k-1} + a_k}{A_{k-1}} \right)^2 - 1 \right) \|z_{k-1}^{x_k}\|_{x_k}^2 \\
& \stackrel{\textcircled{3}}{\geq} \|z_{k-1}^{x_k} - x^*\|_{x_k}^2 + (\xi - 1)\|z_{k-1}^{x_k}\|_{x_k}^2 + \frac{3(\xi - 1)}{2} \left(\frac{1}{1 - \tau_k} - 1 \right) \|z_{k-1}^{x_k}\|_{x_k}^2,
\end{aligned} \tag{8}$$

787 and $\textcircled{2}$ uses the definition of x_k . In $\textcircled{3}$, we used the definition of $\xi = 4\zeta_{2D} - 3$ that implies $\xi - \zeta_{2D} \geq$
788 $\frac{3}{4}(\xi - 1)$ and for $\tau_k \stackrel{\text{def}}{=} a_k/(a_k + A_{k-1})$ we have that $(1 + \frac{a_k}{A_{k-1}})^2 - 1 \geq \frac{2a_k}{A_{k-1}} = 2(\frac{1}{1 - \tau_k} - 1)$.
789 Now, using the second part of Lemma B.1 with $y = y_k$, $x = x_k$ $z^x = -\eta_k v_k^x$, $a^x = z_{k-1}^{x_k}$, so that
790 $z^x + a^x = z_k^{x_k}$ and $z^y + a^y = z_k^{y_k}$ and

$$r = \frac{\|\text{Log}_{x_k}(y_k)\|}{\|z^x\|} = \frac{\lambda\|v_k^x\|}{\eta_k\|v_k^x\|} = \frac{\xi\lambda}{a_k} = \frac{5\xi}{2k + 64\xi} < 5/6 < 1. \tag{9}$$

791 Note that by the choice of parameters and the fact that $r < 1$, the assumptions in [Lemma B.1](#) are
 792 satisfied. Thus, the following holds

$$\|z_k^{x_k} - x^*\|_{x_k}^2 + (\xi - 1)\|z_k^{x_k}\|_{x_k}^2 + \frac{\xi - 1}{2} \left(\frac{r}{1 - r} \right) \|z_{k-1}^{x_k}\|^2 \geq \|z_k^{y_k} - x^*\|_{y_k}^2 + (\xi - 1)\|z_k^{y_k}\|_{y_k}^2. \quad (10)$$

793 Hence, combining (7), (8) and (10) we obtain that it is enough to prove

$$-(1 - \Delta_k) \left(\frac{r}{1 - r} \right) + 3 \left(\frac{1}{1 - \tau_k} - 1 \right) \geq 0,$$

794 The proof will be finished if we prove the result for $\Delta_k = 0$. If we use this last inequality, and the
 795 fact that for $r \leq 5/6$, we have $\frac{r}{1 - r} \leq 3 \left(\frac{1}{1 - 3r/4} - 1 \right)$, we deduce that it suffices to show $\tau_k \geq \frac{3}{4}r$ to
 796 conclude

$$\frac{r}{1 - r} \leq 3 \left(\frac{1}{1 - 3r/4} - 1 \right) \leq 3 \left(\frac{1}{1 - \tau_k} - 1 \right).$$

797 Such inequality, namely $\tau_k \geq \frac{3}{4}r$, is equivalent to $\frac{a_k^2}{\lambda} \geq \frac{3\xi}{4}(a_k + A_{k-1})$ and it holds by [Lemma A.1](#).
 798 \square

799 [Algorithm 1](#) employs a linearly convergent RGD as a subroutine in order to compute [Line 8](#). Below,
 800 we show how this is done and we note that any other linearly convergent algorithm can be used to
 801 solve this step. We first describe a warm start that we will use for RGD. The warm start allows to
 802 know when to stop the subroutine at the same time that it will guarantee fast convergence. One should
 803 think about this lemma as being applied to $h_k(\cdot) \stackrel{\text{def}}{=} f(\cdot) + \frac{1}{2\lambda}d(\cdot, x_k)^2$. Also, note that in that case
 804 we can compute the gradient of h at any point $y \in \mathcal{X}$ as $\nabla h(y) = \nabla f(y) + \frac{1}{\lambda} \text{Log}_y(x_k)$.

805 **Lemma A.4 (Warm start).** *Let \mathcal{M} be a Hadamard manifold, let $x \in \mathcal{M}$, $\mathcal{X} \subset \mathcal{M}$ be a uniquely
 806 geodesic convex set of diameter D and $h : \mathcal{M} \rightarrow \mathbb{R}$ a geodesically convex and L' -smooth function.
 807 Assume access to a projection operator $\mathcal{P}_{\mathcal{X}}$ on \mathcal{X} . Let $x' = \mathcal{P}_{\mathcal{X}}(x)$ and $x^+ \stackrel{\text{def}}{=} \text{Exp}_{x'}(-\frac{1}{L'}\nabla h(x'))$
 808 and $p_0 \stackrel{\text{def}}{=} \mathcal{P}_{\mathcal{X}}(x^+)$ and $D' \stackrel{\text{def}}{=} d(x^+, x') = \|\nabla h(x')\|/L'$. We have that, for all $p \in \mathcal{X}$:*

$$h(p_0) - h(p) \leq \frac{\zeta_{D'}L'}{2}d(x', p)^2 \leq \frac{\zeta_{D'}L'}{2}d(x, p)^2.$$

809 *Proof.* With the notation of the lemma, we have, by smoothness of h , that the following quadratic
 810 $Q : T_{x'}\mathcal{M} \rightarrow \mathbb{R}$, $v \mapsto h(x') + \frac{L'}{2}\|x^+ - v\|_{x'}^2 - \frac{L'}{2}\|x^+ - x'\|_{x'}^2$, induces an upper bound on h in \mathcal{X} ,
 811 via $\text{Exp}_{x'}(\cdot)$. Thus, we have

$$\begin{aligned} -\frac{\zeta_{D'}L'}{2}d(x, p)^2 + h(p_0) &\stackrel{\textcircled{1}}{\leq} -\frac{\zeta_{D'}L'}{2}d(x', p)^2 + h(p_0) \\ &\stackrel{\textcircled{2}}{\leq} -\frac{\zeta_{D'}L'}{2}d(x', p)^2 + Q(\text{Log}_{x'}(p_0)) \\ &\stackrel{\textcircled{3}}{\leq} -\frac{\zeta_{D'}L'}{2}d(x', p)^2 + \left(h(x') + \frac{L'}{2}d(x^+, p_0)^2 - \frac{L'}{2}d(x^+, x')^2 \right) \\ &\stackrel{\textcircled{4}}{\leq} -\frac{\zeta_{D'}L'}{2}d(x', p)^2 + \left(h(x') + \frac{L'}{2}d(x^+, p)^2 - \frac{L'}{2}d(x^+, x')^2 \right) \\ &\stackrel{\textcircled{5}}{\leq} -L'\langle \text{Log}_{x'}(p), \text{Log}_{x'}(x^+) \rangle + h(x') \\ &\stackrel{\textcircled{6}}{=} -L'\langle \text{Log}_{x'}(p), -\frac{1}{L'}\nabla h(x') \rangle + h(x') \\ &\stackrel{\textcircled{7}}{\leq} h(p). \end{aligned}$$

812 We used the projection property of $x' = \mathcal{P}_{\mathcal{X}}(x)$ in $\textcircled{1}$. We used smoothness in $\textcircled{2}$. In $\textcircled{3}$, we
 813 used the first part of [Corollary B.2](#) with $\delta_{D'} = 1$, $r = 1$, $x \leftarrow x'$, $y \leftarrow x^+$, $p \leftarrow p_0$ to bound the

814 estimated distance $\|x^+ - p_0\|_{x'}$ by the actual distance $d(x^+, p_0)$. We used the projection property
815 of $p_0 = \mathcal{P}_{\mathcal{X}}(x^+)$ in (4). In (5), we used the version of [Corollary B.3](#) in [Remark B.4](#). We used the
816 definition of x^+ in (6), and we conclude in (7) by using g-convexity of h . \square

817 Here we finish the computations of the reasoning in [Remark 2.3](#).

818 **Remark A.5.** Let $D'' \stackrel{\text{def}}{=} (L_{f,\mathcal{X}} + 2LD/\zeta_{2D})/L'$, where $L_{f,\mathcal{X}}$ is the Lipschitz constant of
819 f in \mathcal{X} . If we initialize the projected RGD method in [[ZSI16, Theorem 15](#)] with $p_0 \stackrel{\text{def}}{=} \mathcal{P}_{\mathcal{X}}(\text{Exp}_{x'_k}(-\frac{1}{L'}\nabla h_k(x'_k)))$, where $x'_k \stackrel{\text{def}}{=} \mathcal{P}_{\mathcal{X}}(x_k)$, then using (L/ζ_{2D}) -strong g-convexity of h_k to
820 bound $L'd(p_0, y_k^*)^2 \leq 4\zeta_{2D}(h_k(p_0) - h(y_k^*))$, using [Lemma A.4](#) with $x \leftarrow x_k$, $p \leftarrow y_k^*$,
821 $D' \leftarrow \|\nabla h_k(x')\|/L' \leq (\|\nabla f(x')\| + L\|\text{Log}_{x_k}(x')\|/\zeta_{2D})/L' \leq (L_{f,\mathcal{X}} + 2LD/\zeta_{2D})/L' = D''$,

822 and using the guarantees on \mathcal{A} , we have that we find a point y_k satisfying $h_k(y_k) - h_k(y_k^*) \leq$
823 $\frac{\Delta_k d(x_k, y_k^*)^2}{78\lambda}$ in $\tilde{O}(\zeta)$ queries to the gradient and projection oracles. Indeed, the number of queries is
824 given by

$$\begin{aligned} O\left(\zeta_{2D} \log \frac{(h_k(p_0) - h_k(y_k^*)) + L'd(p_0, y_k^*)^2}{\Delta_k d(x_k, y_k^*)^2 / (78\zeta_{2D}/L)}\right) &= O\left(\zeta \log \frac{78\zeta \cdot (1 + 4\zeta_{2D})(\zeta_{D'} L'/2) d(x_k, y_k^*)^2}{L\Delta_k d(x_k, y_k^*)^2}\right) \\ &= O\left(\zeta \log \left(\frac{\zeta \cdot \zeta_{D'}}{\Delta_k}\right)\right) = O\left(\zeta \log \left(\frac{\zeta \cdot \zeta_{D''}}{\Delta_k}\right)\right). \end{aligned}$$

825 Note we know that on the one hand we can stop the algorithm after $O(\zeta \log(\frac{\zeta \cdot \zeta_{D'}}{\Delta_k}))$ iterations which
826 is a value we can compute, including constants, since we can compute D' . On the other hand the
827 worst-case complexity can be expressed as $O(\zeta \log(\frac{\zeta \cdot \zeta_{D''}}{\Delta_k}))$ but we do not need to have access to
828 $L_{f,\mathcal{X}}/L'$. Note that if there is a point $x^* \in \mathcal{X}$ such that $\nabla f(x^*) = 0$, then we have by smoothness
829 that $L_{f,\mathcal{X}} = O(LD)$ and therefore $D'' = O(D)$.

830 Finally, we use [Proposition 2.1](#) and [Remark 2.3](#) to show the final convergence rates for g-convex
831 functions.

832 *Proof (Theorem 2.2).* Given the inequality $(1 - \Delta_k)\psi_k \leq \psi_{k-1}$, proven in [Proposition 2.1](#), we
833 can use ψ_k as a Lyapunov function in order to prove convergence rates of [Algorithm 1](#). It follows
834 straightforwardly by definition of ψ_k , in the following way

$$\begin{aligned} f(y_k) - f(x^*) &\leq \frac{\psi_k}{A_k} \leq \prod_{i=1}^k (1 - \Delta_i)^{-1} \frac{\psi_0}{A_k} \stackrel{\textcircled{1}}{\leq} \frac{2\psi_0}{A_k} \stackrel{\textcircled{2}}{\leq} 2LR_0^2 \left(\frac{A_0}{A_k} + \frac{1}{4LA_k}\right) \\ &= O\left(LR_0^2 \left(\frac{\lambda\xi}{\lambda\left(\frac{k^2+\xi k}{\xi} + \xi\right)} + \frac{1}{\lambda L\left(\frac{k^2+\xi k}{\xi} + \xi\right)}\right)\right) \\ &= O\left(LR_0^2 \left(\frac{\xi^2}{k^2 + \xi k + \xi^2}\right)\right) \stackrel{\textcircled{3}}{=} O\left(\frac{LR_0^2}{k^2} \cdot \zeta^2\right). \end{aligned}$$

835 In (1), we used $\prod_{i=1}^k (1 - \Delta_i) = \prod_{i=1}^k \frac{i(i+2)}{(i+1)^2} = \frac{k+2}{2(k+1)} \geq \frac{1}{2}$. We used smoothness in (2). Note
836 $\frac{\xi-1}{2}\|y_0 - z_0^{y_0}\|_{y_0} = 0$ and $\|z_0^{y_0} - x^*\|_{y_0}^2 = R_0^2$. In (3), we used $\xi = O(\zeta)$ and we dropped some terms
837 in the denominator. Secondly, since the computation of the approximate proximal operator takes
838 $\tilde{O}(\zeta)$ queries to the gradient and projection oracle, cf. [Remark 2.3](#), and $\Delta_k^{-1} \leq \Delta_T^{-1} = (T+1)^2$,
839 then the total number of queries made to these oracles to obtain an ε -minimizer is bounded by
840 $\tilde{O}\left(\zeta^2 \sqrt{\frac{LR_0^2}{\varepsilon}}\right)$. \square

841 We present now the proof that yields an accelerated algorithm for strongly g-convex and smooth
842 functions.

843 *Proof (Theorem 2.4).* The statement of the reduction in [[Mar22, Theorem 7](#)] assumes a function
844 $f: \mathcal{M} \rightarrow \mathbb{R}$ to be optimized has a global minimizer in an unconstrained problem, but the same proof
845 of this theorem works if we have a μ -strongly g-convex and L -smooth function f defined over an open

846 set containing a closed geodesically convex set \mathcal{X} and a minimizer x^* of this function restricted to
847 \mathcal{X} . The reduction provides an algorithm for optimizing f by using $O(\text{Time}_{\text{ns}}(L, \mu, R) \log(\mu R^2/\varepsilon))$
848 queries to the oracle, where $\text{Time}_{\text{ns}}(L, \mu, R)$ is the number of times the oracle is queried by the
849 non-strongly g -convex algorithm if the initial distance is upper bounded by R and if we require
850 accuracy $\mu R^2/4$. In our case, it is $\text{Time}_{\text{ns}}(L, \mu, R) = O(\zeta^2 \log(\zeta^2 \sqrt{L/\mu}) \sqrt{\frac{L}{\mu}}) = O^*(\zeta^2 \sqrt{\frac{L}{\mu}})$, so
851 the result follows. We note that the reverse reduction yields extra geometric penalties but this one
852 does not. \square

853 B Geometric lemmas

854 In this section, we state and prove [Lemma B.5](#), which is used in the proof of [Theorem 2.2](#) to show
855 that the lower bound given by $f(y_k^*) + \langle \tilde{v}_k^y, x - y_k^* \rangle$ that is affine if pulled-back to $T_{y_k^*}$ can be
856 bounded by another function, that is affine if pulled back to x_k . We also include and prove, with
857 some generalizations, some known Riemannian inequalities that are used in Riemannian optimization
858 methods and that we also use. The second part of the following lemma appeared in [\[KY22\]](#). Similarly
859 with the second part of the corollary that follows.

860 In this section, unless otherwise specified, \mathcal{M} is an n -dimensional Riemannian manifold of bounded
861 sectional curvature.

862 **Lemma B.1.** *Let $x, y, p \in \mathcal{M}$ be the vertices of a uniquely geodesic triangle \mathcal{T} of diameter D , and
863 let $z^x \in T_x \mathcal{M}$, $z^y \stackrel{\text{def}}{=} \Gamma_x^y(z^x) + \text{Log}_y(x)$, such that $y = \text{Exp}_x(rz^x)$ for some $r \in [0, 1)$. If we take
864 vectors $a^y \in T_y \mathcal{M}$, $a^x \stackrel{\text{def}}{=} \Gamma_y^x(a^y) \in T_x \mathcal{M}$, then we have the following, for all $\xi \geq \zeta_D$:*

$$\begin{aligned} & \|z^y + a^y - \text{Log}_y(p)\|_y^2 + (\delta_D - 1)\|z^y + a^y\|_y^2 \\ & \geq \|z^x + a^x - \text{Log}_x(p)\|_x^2 + (\delta_D - 1)\|z^x + a^x\|_x^2 - \frac{\xi - \delta_D}{2} \left(\frac{r}{1-r} \right) \|a^x\|_x^2, \end{aligned}$$

865 and

$$\begin{aligned} & \|z^y + a^y - \text{Log}_y(p)\|_y^2 + (\xi - 1)\|z^y + a^y\|_y^2 \\ & \leq \|z^x + a^x - \text{Log}_x(p)\|_x^2 + (\xi - 1)\|z^x + a^x\|_x^2 + \frac{\xi - \delta_D}{2} \left(\frac{r}{1-r} \right) \|a^x\|_x^2. \end{aligned}$$

866 *Proof.* Let γ be the unique geodesic in \mathcal{T} such that $\gamma(0) = x$ and $\gamma(r) = y$. We have $\gamma'(0) = z^x$.
867 Along γ , we define the vector field $V(t) = \Gamma_0^t(\gamma)(z^x - t\gamma'(0))$. Then, it is $V'(t) = -\gamma'(t)$,
868 and $\|V(t)\| = \|a + (1-t)z^x\|$. We will make use of the potential $w : [0, r] \rightarrow \mathbb{R}$ defined as
869 $w(t) = \|\text{Log}_{\gamma(t)}(x) - V(t)\|^2$. We can compute

$$\begin{aligned} \frac{d}{dt} w(t) &= 2\langle D_t(\text{Log}_{\gamma(t)}(x) - V(t)), \text{Log}_{\gamma(t)}(x) - V(t) \rangle \\ &= 2\langle D_t \text{Log}_{\gamma(t)}(x), \text{Log}_{\gamma(t)}(x) \rangle - 2\langle D_t \text{Log}_{\gamma(t)}(x), V(t) \rangle \\ &\quad - 2\langle D_t V(t), \text{Log}_{\gamma(t)}(x) \rangle + 2\langle D_t V(t), V(t) \rangle \\ &= -2\langle D_t(\text{Log}_{\gamma(t)}(x), V(t)) \rangle + 2\langle D_t V(t), V(t) \rangle. \end{aligned} \tag{11}$$

870 Now, we bound the first summand. We use that for the function $\Phi_p(x) = \frac{1}{2}d(x, p)^2$ it holds, for
871 every $\xi \geq \zeta_D$:

$$-\frac{\xi - \delta_D}{2} \|v\|^2 \leq \langle \text{Hess } \Phi_p(\gamma(t))[v] - \frac{\xi + \delta_D}{2} v, v \rangle \leq \frac{\xi - \delta_D}{2} \|v\|^2,$$

872 due to [Fact 1.3](#). So for $\beta \in \{-1, 1\}$ we obtain the following bound:

$$\begin{aligned}
-2\beta \langle D_t \text{Log}_{\gamma(t)}(x), V(t) \rangle &= 2\beta \langle \text{Hess } \Phi_p(\gamma(t))[\gamma'(t)], V(t) \rangle \\
&= 2\beta \langle (\text{Hess } \Phi_p(\gamma(t)) - \frac{\xi + \delta_D}{2} I)[\gamma'(t)], V(t) \rangle + \beta \langle (\xi + \delta_D)\gamma'(t), V(t) \rangle \\
&\leq 2 \|\text{Hess } \Phi_p(\gamma(t)) - \frac{\xi + \delta_D}{2} I\| \cdot \|\gamma'(t)\| \cdot \|V(t)\| + \beta \langle (\xi + \delta_D)\gamma'(t), V(t) \rangle \\
&\leq 2 \frac{\xi - \delta_D}{2} \|\gamma'(t)\| \cdot \|V(t)\| + \beta \langle (\xi + \delta_D)\gamma'(t), V(t) \rangle \\
&\stackrel{\textcircled{1}}{\leq} 2 \frac{\xi - \delta_D}{2} \|z^x\| \cdot \|a + (1-t)z^x\| + \beta \langle (\xi + \delta_D)\gamma'(t), V(t) \rangle
\end{aligned}$$

873 Gauss lemma is used in the last summand of $\textcircled{1}$. Now, if $\beta = -1$, we have

$$\begin{aligned}
-2 \langle D_t \text{Log}_{\gamma(t)}(x), V(t) \rangle &\geq -2 \frac{\xi - \delta_D}{2} \|z^x\| \cdot \|a + (1-t)z^x\| + (\xi + \delta_D) \langle z^x, a + (1-t)z^x \rangle \\
&\stackrel{\textcircled{1}}{\geq} -\frac{\xi - \delta_D}{2(1-t)} (\|(1-t)z^x\|^2 + \|a + (1-t)z^x\|^2) + (\xi - \delta_D) \langle z^x, a + (1-t)z^x \rangle - 2\delta_D \langle -z^x, a + (1-t)b \rangle \\
&\geq -\frac{\xi - \delta_D}{2(1-t)} (\|a\|^2 + 2\langle a + (1-t)z^x \rangle) + (\xi - \delta_D) \langle z^x, a \rangle - 2\delta_D \langle -z^x, a + (1-t)b \rangle \\
&\geq -\frac{\xi - \delta_D}{2(1-t)} \|a\|^2 - 2\delta_D \langle D_t V(t), V(t) \rangle.
\end{aligned} \tag{12}$$

874 On the other hand, analogously, if $\beta = 1$, we have

$$\begin{aligned}
-2 \langle D_t \text{Log}_{\gamma(t)}(x), V(t) \rangle &\leq 2 \frac{\xi - \delta_D}{2} \|z^x\| \cdot \|a + (1-t)z^x\| + (\xi + \delta_D) \langle z^x, a + (1-t)z^x \rangle \\
&\stackrel{\textcircled{1}}{\leq} \frac{\xi - \delta_D}{2(1-t)} (\|(1-t)z^x\|^2 + \|a + (1-t)z^x\|^2) - (\xi - \delta_D) \langle z^x, a + (1-t)z^x \rangle - 2\xi \langle -z^x, a + (1-t)b \rangle \\
&\leq \frac{\xi - \delta_D}{2(1-t)} (\|a\|^2 + 2\langle a + (1-t)z^x \rangle) - (\xi - \delta_D) \langle z^x, a \rangle - 2\xi \langle -z^x, a + (1-t)b \rangle \\
&\leq \frac{\xi - \delta_D}{2(1-t)} \|a\|^2 - 2\xi \langle D_t V(t), V(t) \rangle,
\end{aligned} \tag{13}$$

875 where $\textcircled{1}$ is Young's inequality $2cd \leq c^2 + d^2$. Combining [\(11\)](#), [\(12\)](#), [\(13\)](#), we obtain

$$-\frac{\xi - \delta_D}{2(1-t)} \|a\|^2 - 2(\delta_D - 1) \langle D_t V(t), V(t) \rangle \leq \frac{d}{dt} w(t) \leq \frac{\xi - \delta_D}{2(1-t)} \|a\|^2 - 2(\xi - 1) \langle D_t V(t), V(t) \rangle.$$

876 Integrating between 0 and $r < 1$, it results in

$$\begin{aligned}
\frac{\xi - \delta_D}{2} \log(1-r) \|a\|^2 - (\delta_D - 1) (\|V(r)\|^2 - \|V(0)\|^2) &\leq w(r) - w(0) \\
&\leq -\frac{\xi - \delta_D}{2} \log(1-r) \|a\|^2 - (\xi - 1) (\|V(r)\|^2 - \|V(0)\|^2).
\end{aligned}$$

877 Using the bound $-\log(1-r) \leq \frac{r}{1-r}$ for $r \in [0, 1)$ and using the values of $w(\cdot)$ and $V(\cdot)$, we obtain
878 the result. \square

879 **Corollary B.2.** *Let $x, y, p \in \mathcal{M}$ be the vertices of a uniquely geodesic triangle of diameter D , and
880 let $z^x \in T_x \mathcal{M}$, $z^y \stackrel{\text{def}}{=} \Gamma_x^y(z^x) + \text{Log}_y(x)$, such that $y = \text{Exp}_x(rz^x)$ for some $r \in [0, 1)$. Then, the
881 following holds*

$$\|z^y - \text{Log}_y(p)\|^2 + (\delta_D - 1) \|z^y\|^2 \geq \|z^x - \text{Log}_x(p)\|^2 + (\delta_D - 1) \|z^x\|^2,$$

882 and

$$\|z^y - \text{Log}_y(p)\|^2 + (\zeta_D - 1) \|z^y\|^2 \leq \|z^x - \text{Log}_x(p)\|^2 + (\zeta_D - 1) \|z^x\|^2.$$

883 *Proof.* Use [Lemma B.1](#) with $a^y = 0$. Note that this corollary allows $r = 1$ as well. We obtain this
 884 result, by continuity, by taking a limit when $r \rightarrow 1$. \square

885 The following is a lemma that is already known and is used extensively in Riemannian first-order
 886 optimization. It turns out it is a special case of [Corollary B.2](#).

887 **Corollary B.3 (Cosine-Law Inequalities).** *For the vertices $x, y, p \in \mathcal{M}$ of a uniquely geodesic*
 888 *triangle of diameter D , we have*

$$\langle \text{Log}_x(y), \text{Log}_x(p) \rangle \geq \frac{\delta_D}{2} d(x, y)^2 + \frac{1}{2} d(p, x)^2 - \frac{1}{2} d(p, y)^2.$$

889 *and*

$$\langle \text{Log}_x(y), \text{Log}_x(p) \rangle \leq \frac{\zeta_D}{2} d(x, y)^2 + \frac{1}{2} d(p, x)^2 - \frac{1}{2} d(p, y)^2$$

890 *Proof.* This is [Corollary B.2](#) for $r = 1$. Indeed, given $y \in \mathcal{T}$ we can use [Corollary B.2](#) with $z^x =$
 891 $\text{Log}_x(y)$. Note that in such a case we have $\|z^x\| = d(x, y)$ and $z^y = 0$. Using $\|\text{Log}_y(p)\| = d(y, p)$
 892 *and*

$$\begin{aligned} \|z^x - \text{Log}_x(p)\| &= \|z^x\|^2 - \langle z^x, \text{Log}_x(p) \rangle + \|\text{Log}_x(p)\|^2 \\ &= d(x, y)^2 - 2\langle \text{Log}_x(y), \text{Log}_x(p) \rangle + d(p, x)^2, \end{aligned}$$

893 we obtain the result. \square

894 **Remark B.4.** *Actually, in Hadamard manifolds, if we substitute the constants δ_D and ζ_D in the*
 895 *previous [Corollary B.3](#) by the tighter constants $\delta_{d(p,x)}$ and $\zeta_{d(p,x)}$, the result also holds. See [\[ZS16\]](#).*

896 We now proceed to prove a lemma that intuitively says that solving the exact proximal point problem
 897 can be used to lower bound f . One should think about the following lemma as being applied
 898 to $y \leftarrow y_k^*$, $x \leftarrow x_k$. Compare the result of the following lemma with the Euclidean equality
 899 $\langle g, p - y \rangle = \langle g, p - x \rangle + \|g\|^2$, for $g = x - y$ and $x, y, p \in \mathbb{R}^n$.

900 **Lemma B.5.** *Let $x, y, p \in \mathcal{M}$ be the vertices of a uniquely geodesic triangle \mathcal{T} of diameter D . Define*
 901 *the vectors $g \stackrel{\text{def}}{=} \text{Log}_y(x) \in T_y\mathcal{M}$ and $g^x = \Gamma_y^x(g) = -\text{Log}_x(y) \in T_x\mathcal{M}$. Then we have*

$$\langle g, \text{Log}_y(p) \rangle \geq \langle g^x, \text{Log}_x(p) \rangle + \delta_D \|g\|^2,$$

902 *and*

$$\langle g, \text{Log}_y(p) \rangle \leq \langle g^x, \text{Log}_x(p) \rangle + \zeta_D \|g\|^2.$$

903 *Proof* ([Lemma B.5](#)). Using the definition of g , we have ① below, by the first part of [Corollary B.3](#):

$$\begin{aligned} \langle g, \text{Log}_y(p) \rangle &\stackrel{\textcircled{1}}{\geq} \frac{\delta_D}{2} \|g\|^2 + \frac{d(y, p)^2}{2} - \frac{d(x, p)^2}{2} \\ &\stackrel{\textcircled{2}}{\geq} \langle g^x, \text{Log}_x(p) \rangle + \delta_D \|g^x\|^2, \end{aligned}$$

904 *and in* ② *we used [Corollary B.3](#) again but with a different choice of vertices so we have* $\frac{d(y, p)^2}{2} \geq$
 905 $\frac{\delta_D}{2} \|g^x\|^2 + \frac{d(x, p)^2}{2} + \langle g^x, \text{Log}_x(p) \rangle$.

906 The proof of the second part is analogous: using the definition of g , we have ① below, by the second
 907 part of [Corollary B.3](#):

$$\begin{aligned} \langle g, \text{Log}_y(p) \rangle &\stackrel{\textcircled{1}}{\leq} \frac{\zeta_D}{2} \|g\|^2 + \frac{d(y, p)^2}{2} - \frac{d(x, p)^2}{2} \\ &\stackrel{\textcircled{2}}{\leq} \langle g^x, \text{Log}_x(p) \rangle + \zeta_D \|g^x\|^2, \end{aligned}$$

908 *and in* ② *we used [Corollary B.3](#) again but with a different choice of vertices so we have* $\frac{d(y, p)^2}{2} \leq$
 909 $\frac{\zeta_D}{2} \|g^x\|^2 + \frac{d(x, p)^2}{2} + \langle g^x, \text{Log}_x(p) \rangle$. \square

910 **C Other subroutines**

911 We provide two other subroutines that optimize functions that are μ -strongly g -convex and L -smooth
 912 with linear rates and thus they can be used as subroutines for Line 8 in Algorithm 1. This yields
 913 accelerated algorithms for each of them.

914 For the first subroutine, we change the analysis but use the same algorithm as ZS16: Projected
 915 Riemannian Gradient descent $x_{t+1} \leftarrow P_X(\text{Exp}_{x_t}(-\eta\nabla f(x_t)))$ but we set learning rate $\eta \stackrel{\text{def}}{=} (2 -$
 916 $\zeta_D)/L$. Let $\tilde{x}_{t+1} \stackrel{\text{def}}{=} \text{Exp}_{x_t}(-\eta\nabla f(x_t))$. First we show the following inequality that results from
 917 applying smoothness to the first part and strong g -convexity to the second one.

$$\begin{aligned}
 0 &\leq f(\tilde{x}_{t+1}) - f(x^*) = f(\tilde{x}_{t+1}) - f(x_t) + f(x_t) - f(x^*) \\
 &\leq \langle \nabla f(x_t), \tilde{x}_{t+1} - x_t \rangle + \frac{L}{2} \|\tilde{x}_{t+1} - x_t\|_{x_t}^2 + \langle \nabla f(x_t), x_t - x^* \rangle - \frac{\mu}{2} \|x_t - x^*\|_{x_t}^2 \\
 &= \langle \nabla f(x_t), \tilde{x}_{t+1} - x^* \rangle + \frac{L\eta^2}{2} \|\nabla f(x_t)\|^2 - \frac{\mu}{2} \|x_t - x^*\|_{x_t}^2 \\
 &= \langle \nabla f(x_t), x_t - x^* \rangle + \left(\frac{L\eta^2}{2} - \eta\right) \|\nabla f(x_t)\|^2 - \frac{\mu}{2} \|x_t - x^*\|_{x_t}^2.
 \end{aligned} \tag{14}$$

918 Now, we have the following bound, bounding the distance to the minimizer, from which we will
 919 derive convergence rates for projected RGD:

$$\begin{aligned}
 d(\tilde{x}_{t+1}, x^*)^2 &\stackrel{\textcircled{1}}{\leq} (\zeta - 1)\eta^2 \|\nabla f(x_t)\|^2 + \|x^* - \tilde{x}_{t+1}\|_{x_t}^2 \\
 &\stackrel{\textcircled{2}}{\leq} \|x^* - x_t\|_{x_t}^2 + 2\eta \langle \nabla f(x_t), x^* - x_t \rangle + \zeta\eta^2 \|\nabla f(x_t)\|^2 \\
 &\stackrel{\textcircled{3}}{\leq} \left(2\eta - \frac{\zeta\eta}{1 - \frac{L\eta}{2}}\right) \langle \nabla f(x_t), x^* - x_t \rangle + \left(1 - \frac{\mu\zeta\eta}{1 - \frac{L\eta}{2}}\right) \|x^* - x_t\|_{x_t}^2.
 \end{aligned} \tag{15}$$

920 where in $\textcircled{1}$ we used the Euclidean cosine theorem along with Corollary B.3. Inequality $\textcircled{2}$ develops
 921 the square $\|x^* - \tilde{x}_{t+1}\|_{x_t}^2 = \|x^* - x_t - \eta\nabla f(x_t)\|_{x_t}^2$ and $\textcircled{3}$ uses (14), where the inequality has
 922 been multiplied by $-\zeta\eta^2(L\eta^2/2 - \eta)^{-1} = \frac{\zeta\eta}{1 - \frac{L\eta}{2}}$ (≥ 0 , since we assume $\eta \in [0, 2/L]$) in both sides.

923 Now, since $\langle \nabla f(x_t), x^* - x_t \rangle \leq 0$, we want to make the factor alongside it be ≥ 0 in order to drop
 924 it. That means, it should be $2\eta - \frac{\zeta\eta}{1 - \frac{L\eta}{2}} \geq 0$ which is equivalent to $\eta \leq \frac{2-\zeta}{L}$. By setting η exactly to
 925 the value $\frac{2-\zeta}{L}$ and assuming $\zeta < 2$, we have $\frac{\zeta\eta}{1 - \frac{L\eta}{2}} = 2(2 - \zeta)/L$ and so we can conclude:

$$d(x_{t+1}, x^*)^2 \leq d(\tilde{x}_{t+1}, x^*)^2 \leq \left(1 - \frac{2\mu(2 - \zeta)}{L}\right) d(x_t, x^*)^2.$$

926 which is linear convergence, as desired.

927 For the second subroutine, we assume access to the operation

$$x_{t+1} = \arg \min_{y \in \mathcal{X}} \{ \langle \nabla f(x_t), y - x_t \rangle_{x_t} + \frac{L}{2} d(x_t, y)^2 \},$$

928 and define the algorithm as the sequential application of it. This subproblem, in the Euclidean case,
 929 is equivalent to the projection operator of $\tilde{x}_{t+1} = \text{Exp}_{x_t}(-\eta\nabla f(x_t))$. However, in the Riemannian
 930 case, this and the metric-projection operator $P_{\mathcal{X}}(x_{t+1})$ are two different things in general. Define the
 931 notation $\phi(x) \stackrel{\text{def}}{=} (f + I_{\mathcal{X}})(x)$. Then, we have

$$\begin{aligned}
\phi(x_{t+1}) &\stackrel{\textcircled{1}}{\leq} m_L(x_t, x_{t+1}) \\
&= \min_{x \in \mathcal{M}} \left\{ f(x_t) + \langle \nabla f(x_t), x - x_t \rangle_{x_t} + \frac{L}{2} d(x, x_t)^2 + I_{\mathcal{X}}(x) \right\} \\
&\stackrel{\textcircled{2}}{\leq} \min_{x \in \mathcal{M}} \left\{ f(x) + \frac{L}{2} d(x, x_t)^2 + I_{\mathcal{X}}(x) \right\} \\
&= \min_{x \in \mathcal{M}} \left\{ \phi(x) + \frac{L}{2} d(x, x_t)^2 \right\} \\
&\stackrel{\textcircled{3}}{\leq} \min_{\alpha \in [0,1]} \left\{ \alpha \phi(x^*) + (1 - \alpha) \phi(x_t) + \frac{L\alpha^2}{2} d(x^*, x_t)^2 \right\} \\
&\stackrel{\textcircled{4}}{\leq} \min_{\alpha \in [0,1]} \left\{ \phi(x_t) - \alpha \left(1 - \alpha \frac{L}{\mu} \right) (\phi(x_t) - \phi(x^*)) \right\} \\
&\stackrel{\textcircled{5}}{=} \phi(x_t) - \frac{\mu}{2L} (\phi(x_t) - \phi(x^*)).
\end{aligned}$$

932 Above, $\textcircled{1}$ holds by smoothness and $\textcircled{2}$ holds by g-convexity of f (I thought maybe using strong
933 convexity one can improve but it is not by much, it results in convergence rates of $O((\frac{L}{\mu} - 1) \log(1/\varepsilon))$
934 instead of $O(\frac{L}{\mu} \log(1/\varepsilon))$. So I am not using it). Inequality $\textcircled{3}$ results from restricting the minimum to
935 the geodesic segment between x^* and x_t and uses g-convexity of ψ . In $\textcircled{4}$, we used strong convexity
936 of ϕ to bound $\frac{\mu}{2} d(x^*, y_k)^2 \leq \phi(x_t) - \phi(x^*)$. Finally, in $\textcircled{5}$ we substituted α by the value that
937 minimizes the expression, which is $\mu/2L$.

938 Subtracting $\phi(x^*)$ to the inequality above, we obtain $\phi(x_{t+1}) - \phi(x^*) \leq (1 - \frac{\mu}{2L}) (\phi(x_t) - \phi(x^*))$.
939 As we wanted to prove, there is linear convergence.