

1 **A VARIATIONAL FORMULATION OF**
2 **ACCELERATED OPTIMIZATION ON RIEMANNIAN MANIFOLDS**

3 VALENTIN DURUISSEAUX AND MELVIN LEOK

4 **ABSTRACT.** It was shown recently by Su et al. [22] that Nesterov’s accelerated gradient method for
5 minimizing a smooth convex function f can be thought of as the time discretization of a second-
6 order ODE, and that $f(x(t))$ converges to its optimal value at a rate of $\mathcal{O}(1/t^2)$ along any trajectory
7 $x(t)$ of this ODE. A variational formulation was introduced in Wibisono et al. [24] which allowed
8 for accelerated convergence at a rate of $\mathcal{O}(1/t^p)$, for arbitrary $p > 0$, in normed vector spaces.
9 This framework was exploited in Duruisseaux et al. [7] using time-adaptive geometric integrators
10 to design efficient explicit algorithms for symplectic accelerated optimization. In Alimisis et al.
11 [4], a second-order ODE was proposed as the continuous-time limit of a Riemannian accelerated
12 algorithm, and it was shown that the objective function $f(x(t))$ converges to its optimal value at
13 a rate of $\mathcal{O}(1/t^2)$ along solutions of this ODE, thereby generalizing the earlier Euclidean result
14 to the Riemannian manifold setting. In this paper, we show that on Riemannian manifolds, the
15 convergence rate of $f(x(t))$ to its optimal value can also be accelerated to an arbitrary convergence
16 rate $\mathcal{O}(1/t^p)$, by considering a family of time-dependent Bregman Lagrangian and Hamiltonian
17 systems on Riemannian manifolds. This generalizes the results of [24] to Riemannian manifolds
18 and also provides a variational framework for accelerated optimization on Riemannian manifolds.
19 In particular, we will establish results for objective functions on Riemannian manifolds that are
20 geodesically convex, weakly-quasi-convex, and strongly convex. An approach based on the time-
21 invariance property of the family of Bregman Lagrangians and Hamiltonians was used to construct
22 very efficient optimization algorithms in [7], and we establish a similar time-invariance property in
23 the Riemannian setting. This lays the foundation for constructing similarly efficient optimization
24 algorithms on Riemannian manifolds, once the Riemannian analogue of time-adaptive Hamiltonian
25 variational integrators has been developed. The experience with the numerical discretization of
26 variational accelerated optimization flows on vector spaces suggests that the combination of time-
27 adaptivity and symplecticity is important for the efficient, robust, and stable discretization of these
28 variational flows describing accelerated optimization. One expects that a geometric numerical
29 integrator that is time-adaptive, symplectic, and Riemannian manifold preserving will yield a class
30 of similarly promising optimization algorithms on manifolds.

31 1. INTRODUCTION

32 Efficient optimization has become one of the major concerns in data analysis. Many machine
33 learning algorithms are designed around the minimization of a loss function or the maximization
34 of a likelihood function. Due to the ever-growing scale of the data sets and size of the problems,
35 there has been a lot of focus on first-order optimization algorithms because of their low cost per
36 iteration. The first gradient descent algorithm was proposed in [6] by Cauchy to deal with the very
37 large systems of equations he was facing when trying to simulate orbits of celestial bodies, and
38 many gradient-based optimization methods have been proposed since Cauchy’s work in 1847.

39 In 1983, Nesterov’s accelerated gradient method was introduced in [18], and was shown to con-
40 verge in $\mathcal{O}(1/k^2)$ to the minimum of the convex objective function f , improving on the $\mathcal{O}(1/k)$
41 convergence rate exhibited by the standard gradient descent methods. This $\mathcal{O}(1/k^2)$ convergence
42 rate was shown in [19] to be optimal among first-order methods using only information about ∇f
43 at consecutive iterates. This phenomenon in which an algorithm displays this improved rate of
44 convergence is referred to as acceleration, and other accelerated algorithms have been derived since
45 Nesterov’s algorithm, such as accelerated mirror descent [17] and accelerated cubic-regularized
46 Newton’s method [20]. More recently, it was shown in [22] that Nesterov’s accelerated gradient

method limits to a second order ODE, as the timestep goes to 0, and that the objective function $f(x(t))$ converges to its optimal value at a rate of $\mathcal{O}(1/t^2)$ along the trajectories of this ODE. It was then shown in [24] that in continuous time, the convergence rate of $f(x(t))$ can be accelerated to an arbitrary convergence rate $\mathcal{O}(1/t^p)$ in normed spaces, by considering flow maps generated by a family of time-dependent Bregman Lagrangian and Hamiltonian systems which is closed under time rescaling. This variational framework and the time-invariance property of the family of Bregman Lagrangians was then exploited in [7] using time-adaptive geometric integrators to design efficient explicit algorithms for symplectic accelerated optimization. It was observed that a careful use of adaptivity and symplecticity could result in a significant gain in computational efficiency.

In the past few years, there has been some effort to derive accelerated optimization algorithms in the Riemannian manifold setting [2–4; 14; 25; 26]. In [4], a second order ODE was proposed as the continuous-time limit of a Riemannian accelerated algorithm, and it was shown that the objective function $f(x(t))$ converges to its optimal value at a rate of $\mathcal{O}(1/t^2)$ along solutions of this ODE, generalizing the Euclidean result obtained in [22] to the Riemannian manifold setting.

In this paper, we show that in continuous time, the convergence rate of $f(x(t))$ to its optimal value can be accelerated to an arbitrary convergence rate $\mathcal{O}(1/t^p)$ on Riemannian manifolds, thereby generalizing the results of [24] to the Riemannian setting. This is achieved by considering a family of time-dependent Bregman Lagrangian and Hamiltonian systems on Riemannian manifolds. This also provides a variational framework for accelerated optimization on Riemannian manifolds, generalizing the normed vector space variational formulation of accelerated optimization introduced in [24]. We will then illustrate the derived theoretical convergence rates by integrating the Bregman Euler–Lagrange equations using a simple numerical scheme to solve eigenvalue and distance minimization problems on Riemannian manifolds. Finally, we will show that the family of Bregman dynamics is closed under time rescaling, and we will draw inspiration from the approach introduced in [7] to take advantage of this invariance property via a carefully chosen Poincaré transformation that will allow for the integration of higher-order Bregman dynamics while benefiting from the computational efficiency of integrating lower-order Bregman dynamics on Riemannian manifolds.

2. DEFINITIONS AND PRELIMINARIES

We first introduce the main notions from Riemannian geometry and Lagrangian and Hamiltonian mechanics that will be used throughout this paper (see [4; 8; 9; 11; 12; 15] for more details).

2.1. Riemannian Geometry.

Definition 2.1. *Suppose we have a Riemannian manifold \mathcal{Q} with Riemannian metric $g(\cdot, \cdot) = \langle \cdot, \cdot \rangle$, represented by the positive-definite symmetric matrix (g_{ij}) in local coordinates. Then, we define the **musical isomorphism** $g^\flat : T\mathcal{Q} \rightarrow T^*\mathcal{Q}$ via*

$$g^\flat(u)(v) = g_p(u, v) \quad \forall p \in \mathcal{Q} \text{ and } \forall u, v \in T_p\mathcal{Q},$$

and its **inverse musical isomorphism** $g^\sharp : T^*\mathcal{Q} \rightarrow T\mathcal{Q}$. The Riemannian metric $g(\cdot, \cdot) = \langle \cdot, \cdot \rangle$ induces a **fiber metric** $g^*(\cdot, \cdot) = \langle\langle \cdot, \cdot \rangle\rangle$ on $T^*\mathcal{Q}$ via

$$\langle\langle u, v \rangle\rangle = \langle g^\sharp(u), g^\sharp(v) \rangle \quad \forall u, v \in T^*\mathcal{Q},$$

represented by the positive definite symmetric matrix (g^{ij}) in local coordinates, which is the inverse of the Riemannian metric matrix (g_{ij}) .

Definition 2.2. *The **Riemannian gradient** $\text{grad}f(q) \in T_q\mathcal{Q}$ at a point $q \in \mathcal{Q}$ of a smooth function $f : \mathcal{Q} \rightarrow \mathbb{R}$ is the tangent vector at q such that*

$$\langle \text{grad}f(q), u \rangle = df(q)u \quad \forall u \in T_q\mathcal{Q},$$

where df is the differential of f .

91 **Definition 2.3.** A **vector field** on a Riemannian manifold \mathcal{Q} is a map $X : \mathcal{Q} \rightarrow T\mathcal{Q}$ such that
92 $X(q) \in T_q\mathcal{Q}$ for all $q \in \mathcal{Q}$. The set of all vector fields on \mathcal{Q} is denoted $\mathcal{X}(\mathcal{Q})$. The **integral curve**
93 at q of $X \in \mathcal{X}(\mathcal{Q})$ is the smooth curve c on \mathcal{Q} such that $c(0) = q$ and $c'(t) = X(c(t))$.

94 **Definition 2.4.** A **geodesic** in a Riemannian manifold \mathcal{Q} is a parametrized curve $\gamma : [0, 1] \rightarrow \mathcal{Q}$
95 which is of minimal local length. It can be thought of as a curve having zero “acceleration” or
96 constant “speed”, that is as a generalization of the notion of straight line from Euclidean spaces
97 to Riemannian manifolds. Given two points $q, \tilde{q} \in \mathcal{Q}$, a vector in $T_q\mathcal{Q}$ can be transported to $T_{\tilde{q}}\mathcal{Q}$
98 along a geodesic γ by an operation $\Gamma(\gamma)_{\tilde{q}}^q : T_q\mathcal{Q} \rightarrow T_{\tilde{q}}\mathcal{Q}$ called **parallel transport along** γ . We
99 will simply write $\Gamma_{\tilde{q}}^q$ to denote the parallel transport along some geodesic connecting the two points
100 $q, \tilde{q} \in \mathcal{Q}$, and given $A \in \mathcal{X}(\mathcal{Q})$, we will denote by $\Gamma(A)$ the parallel transport along integral curves
101 of A . Note that parallel transport preserves inner products: given a geodesic γ from $q \in \mathcal{Q}$ to $\tilde{q} \in \mathcal{Q}$,

$$102 \quad g_q(u, v) = g_{\tilde{q}}(\Gamma(\gamma)_{\tilde{q}}^q u, \Gamma(\gamma)_{\tilde{q}}^q v) \quad \forall u, v \in T_q\mathcal{Q}.$$

103 **Definition 2.5.** Given $X, Y \in \mathcal{X}(\mathcal{Q})$, the **covariant derivative** $\nabla_X Y \in \mathcal{X}(\mathcal{Q})$ of Y along X is

$$104 \quad \nabla_X Y(q) = \lim_{h \rightarrow 0} \frac{\Gamma(\gamma)_{\gamma(h)}^q Y(\gamma(h)) - Y(q)}{h},$$

105 where γ is the unique integral curve of X such that $\gamma(0) = q$, for any $q \in \mathcal{Q}$.

106 **Definition 2.6.** A function $f : \mathcal{Q} \rightarrow \mathbb{R}$ is called **L -smooth** if for any two points $q, \tilde{q} \in \mathcal{Q}$ and
107 geodesic γ connecting them,

$$108 \quad \|\text{grad}f(q) - \Gamma(\gamma)_{\tilde{q}}^q \text{grad}f(\tilde{q})\| \leq L \text{ length}(\gamma).$$

109 **Definition 2.7.** The **Riemannian Exponential map** $\text{Exp}_q : T_q\mathcal{Q} \rightarrow \mathcal{Q}$ at $q \in \mathcal{Q}$ is defined via

$$110 \quad \text{Exp}_q(v) = \gamma_v(1),$$

111 where γ_v is the unique geodesic in \mathcal{Q} such that $\gamma_v(0) = q$ and $\gamma_v'(0) = v$, for any $v \in T_q\mathcal{Q}$.

112 Exp_q is a diffeomorphism in some neighborhood $U \subset T_q\mathcal{Q}$ containing 0, so we can define its inverse
113 map, the **Riemannian Logarithm map** $\text{Log}_p : \text{Exp}_q(U) \rightarrow T_q\mathcal{Q}$.

114 **Definition 2.8.** Given a Riemannian manifold \mathcal{Q} with sectional curvature bounded below by K_{\min} ,
115 and an upper bound D for the diameter of the considered domain, define

$$116 \quad \zeta = \begin{cases} \sqrt{-K_{\min}} D \coth(\sqrt{-K_{\min}} D) & \text{if } K_{\min} < 0 \\ 1 & \text{if } K_{\min} \geq 0 \end{cases}. \quad (2.1)$$

117 Note that $\zeta \geq 1$ since $x \coth x \geq 1$ for all real values of x .

118 2.2. Convexity in Riemannian Manifolds.

119 **Definition 2.9.** A subset A of a Riemannian manifold \mathcal{Q} is called **geodesically uniquely convex**
120 if every two points of A are connected by a unique geodesic in A . A function $f : \mathcal{Q} \rightarrow \mathbb{R}$ is called
121 **geodesically convex** if for any two points $q, \tilde{q} \in \mathcal{Q}$ and geodesic γ connecting them,

$$122 \quad f(\gamma(t)) \leq (1-t)f(q) + tf(\tilde{q}) \quad \forall t \in [0, 1].$$

123 Note that if f is a smooth geodesically convex function on a geodesically uniquely convex subset A
124 of a Riemannian manifold, then

$$125 \quad f(q) - f(\tilde{q}) \geq \langle \text{grad}f(\tilde{q}), \text{Log}_{\tilde{q}}(q) \rangle \quad \forall q, \tilde{q} \in A.$$

126 A function $f : A \rightarrow \mathbb{R}$ is called **geodesically α -weakly-quasi-convex** (α -WQC) with respect to
127 $q \in \mathcal{Q}$ for some $\alpha \in (0, 1]$ if

$$128 \quad \alpha(f(q) - f(\tilde{q})) \geq \langle \text{grad}f(\tilde{q}), \text{Log}_{\tilde{q}}(q) \rangle \quad \forall \tilde{q} \in A.$$

129 A function $f : A \rightarrow \mathbb{R}$ is called **geodesically μ -strongly-convex** (μ -SC) for some $\mu > 0$ if

$$130 \quad f(q) - f(\tilde{q}) \geq \langle \text{grad}f(\tilde{q}), \text{Log}_{\tilde{q}}(q) \rangle + \frac{\mu}{2} \|\text{Log}_{\tilde{q}}(q)\|^2 \quad \forall q, \tilde{q} \in A.$$

131 A local minimum of a geodesically convex or α -WQC function is also a global minimum, and a
132 geodesically strongly convex function either has no minimum or a unique global minimum.

133 **2.3. Lagrangian and Hamiltonian Mechanics.** Given a n -dimensional Riemannian manifold
134 \mathcal{Q} with local coordinates (q^1, \dots, q^n) , a **Lagrangian** is a function $L : T\mathcal{Q} \times \mathbb{R} \rightarrow \mathbb{R}$. The **action**
135 **integral** \mathcal{S} is defined to be the functional

$$136 \quad \mathcal{S}(q) = \int_0^T L(q, \dot{q}, t) dt, \quad (2.2)$$

137 over the space of smooth curves $q : [0, T] \rightarrow \mathcal{Q}$. **Hamilton's Variational Principle** states that
138 $\delta\mathcal{S} = 0$ where the variation $\delta\mathcal{S}$ is induced by an infinitesimal variation δq of the trajectory q that
139 vanishes at the endpoints. Hamilton's Variational Principle can be shown to be equivalent to the
140 **Euler–Lagrange equations**

$$141 \quad \frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}^k} \right) = \frac{\partial L}{\partial q^k} \quad \text{for } k = 1, \dots, n. \quad (2.3)$$

142 The **Legendre transform** $\mathbb{F}L : T\mathcal{Q} \rightarrow T^*\mathcal{Q}$ of L is defined fiberwise via $\mathbb{F}L : (q^i, \dot{q}^i) \mapsto (q^i, p_i)$
143 where $p_i = \frac{\partial L}{\partial \dot{q}^i} \in T^*\mathcal{Q}$ is the **conjugate momentum** of q^i . We can then define the associated
144 **Hamiltonian** $H : T^*\mathcal{Q} \rightarrow \mathbb{R}$ via

$$145 \quad H(q, p, t) = \sum_{j=1}^n p_j \dot{q}^j - L(q, \dot{q}, t) \Big|_{p_i = \frac{\partial L}{\partial \dot{q}^i}}. \quad (2.4)$$

146 We can also define a Hamiltonian Variational Principle on the Hamiltonian side in momentum
147 phase space

$$148 \quad \delta \int_0^T \sum_{j=1}^n [p_j \dot{q}^j - H(q, p, t)] dt = 0,$$

149 where the variation is induced by an infinitesimal variation δq of the trajectory q that vanishes at
150 the endpoints. This is equivalent to **Hamilton's equations**, given by

$$151 \quad \dot{p}_k = -\frac{\partial H}{\partial q^k}(p, q), \quad \dot{q}^k = \frac{\partial H}{\partial p_k}(p, q) \quad \text{for } k = 1, \dots, n, \quad (2.5)$$

152 which can also be shown to be equivalent to the Euler–Lagrange equations (2.3).

153 3. VARIATIONAL FORMULATION AND CONVERGENCE RATES

154 Throughout this paper, we will make the following assumptions on the function $f : \mathcal{Q} \rightarrow \mathbb{R}$ to
155 be minimized and on the ambient Riemannian manifold \mathcal{Q} , which are standard assumptions in
156 Riemannian optimization [3; 4; 25; 26]:

157 **Assumption 1.** *Solutions of the differential equations derived in this paper remain inside a geodesi-*
158 *cally uniquely convex subset A of a complete Riemannian manifold \mathcal{Q} (i.e. any two points in \mathcal{Q}*
159 *can be connected by a geodesic), such that $\text{diam}(A)$ is bounded above by some constant D , that*
160 *the sectional curvature is bounded from below by K_{\min} on A , and that Exp_q is well-defined for any*
161 *$q \in A$, and its inverse Log_q is well-defined and differentiable on A for any $q \in A$. Furthermore, f is*
162 *bounded below, geodesically L -smooth and all its minima are inside A .*

163 **3.1. Convex Case.** Suppose $f : \mathcal{Q} \rightarrow \mathbb{R}$ is a geodesically convex function, and that Assumption 1
164 holds true. We define the corresponding p -Bregman Lagrangian $\mathcal{L}_p^C : T\mathcal{Q} \times \mathbb{R} \rightarrow \mathbb{R}$ for $p > 0$ via

$$165 \quad \boxed{\mathcal{L}_p^C(X, V, t) = \frac{t^{\zeta p+1}}{2p} \langle V, V \rangle - Cpt^{(\zeta+1)p-1} f(X)}, \quad (3.1)$$

166 and the corresponding p -Bregman Hamiltonian $\mathcal{H}_p^C : T^*\mathcal{Q} \times \mathbb{R} \rightarrow \mathbb{R}$ is given by

$$167 \quad \boxed{\mathcal{H}_p^C(X, R, t) = \frac{p}{2t^{\zeta p+1}} \langle R, R \rangle + Cpt^{(\zeta+1)p-1} f(X)}, \quad (3.2)$$

168 where $X \in \mathcal{Q}$ denotes position on the manifold \mathcal{Q} , V and R are velocity vector and momentum
169 covector fields, t is the time variable, C is a constant, and ζ is given by equation (2.1). This
170 choice of Bregman Lagrangian is inspired by the results of [5; 7; 24], and can be thought of as a
171 generalization of the normed space p -Bregman Lagrangians and Hamiltonians

$$172 \quad L(X, V, t) = \frac{t^{p+1}}{2p} \langle V, V \rangle - Cpt^{2p-1} f(X), \quad H(X, R, t) = \frac{p}{2t^{p+1}} \langle R, R \rangle + Cpt^{2p-1} f(X), \quad (3.3)$$

173 obtained in [7], where the structure of the Riemannian manifold \mathcal{Q} has now been incorporated
174 through the constant ζ . These p -Bregman Lagrangians and Hamiltonians arise from the Bregman
175 Lagrangians and Hamiltonians introduced in [24],

$$176 \quad \mathcal{L}_{\alpha, \beta, \gamma}(x, v, t) = e^{\alpha(t)+\gamma(t)} \left[D_h(x + e^{-\alpha(t)}v, x) - e^{\beta(t)} f(x) \right], \quad (3.4)$$

$$177 \quad \mathcal{H}_{\alpha, \beta, \gamma}(x, r, t) = e^{\alpha(t)+\gamma(t)} \left[D_{h^*}(\nabla h(x) + e^{-\gamma(t)}r, \nabla h(x)) + e^{\beta(t)} f(x) \right], \quad (3.5)$$

179 where the Bregman divergence D_h is given by

$$180 \quad D_h(x, y) = h(y) - h(x) - \langle \nabla h(x), y - x \rangle,$$

181 for the chosen convex, continuously differentiable function $h(x) = \frac{1}{2} \langle x, x \rangle$, and where the Legendre
182 transform (or convex dual function) h^* is given by $h^* = \sup_{v \in T\mathcal{X}} [\langle r, v \rangle - h(v)]$. The parameter
183 functions α, β, γ are chosen to be

$$184 \quad \alpha_t = \log \zeta p - \log t, \quad \beta_t = p \log t + \log C - 2 \log \zeta, \quad \gamma_t = \zeta p \log t + \log \zeta, \quad (3.6)$$

185 and these satisfy the ideal scaling conditions $\dot{\beta}_t \leq e^{\alpha_t}$ and $\dot{\gamma}_t = e^{\alpha_t}$. The ideal scaling conditions were
186 necessary conditions introduced in [24] for the Bregman Lagrangians and Hamiltonians to have flows
187 that converge to the minimizer at the rate $\mathcal{O}(e^{-\beta_t})$. In this paper, we will simplify the exposition
188 by focusing on the more practically relevant case of Bregman Lagrangians and Hamiltonians that
189 are parametrized by $p > 0$ which achieve a convergence rate of $\mathcal{O}(1/t^p)$.

190 **Theorem 3.1.** *The p -Bregman Euler–Lagrange equation corresponding to \mathcal{L}_p^C is given by*

$$191 \quad \boxed{\nabla_{\dot{X}} \dot{X} + \frac{\zeta p + 1}{t} \dot{X} + Cp^2 t^{p-2} \text{grad}f(X) = 0}. \quad (3.7)$$

192 *Proof.* See Appendix A.1, with $\lambda = \zeta$.

193 **Theorem 3.2.** *Suppose $f : \mathcal{Q} \rightarrow \mathbb{R}$ is a geodesically convex function, and Assumption 1 is satisfied.*
194 *Then, the p -Bregman Euler–Lagrange equation (3.7) has a solution, and any solution converges to*
195 *a minimizer x^* of f with rate*

$$196 \quad \boxed{f(X(t)) - f(x^*) \leq \frac{\zeta \|\text{Log}_{x_0}(x^*)\|^2}{2Ct^p}}. \quad (3.8)$$

197 *Proof.* See Appendix B.1 for the existence of a solution and Appendix C.1 for the convergence rate.

198 Note that this theorem reduces to Theorem 5 from [4] when $p = 2$ and $C = 1/4$.

199 **3.2. Weakly-Quasi-Convex Case.** Suppose $f : \mathcal{Q} \rightarrow \mathbb{R}$ is a geodesically α -weakly-quasi-convex
 200 function, and suppose that Assumption 1 is satisfied. We define the corresponding p -Bregman
 201 Lagrangian $\mathcal{L}_p^{WQC} : T\mathcal{Q} \times \mathbb{R} \rightarrow \mathbb{R}$ for $p > 0$ via

$$202 \quad \mathcal{L}_p^{WQC}(X, V, t) = \frac{t^{\frac{\zeta}{\alpha}p+1}}{2p} \langle V, V \rangle - Cpt^{\left(\frac{\zeta}{\alpha}+1\right)p-1} f(X), \quad (3.9)$$

203 and the corresponding p -Bregman Hamiltonian $\mathcal{H}_p^{WQC} : T^*\mathcal{Q} \times \mathbb{R} \rightarrow \mathbb{R}$ is given by

$$204 \quad \mathcal{H}_p^{WQC}(X, R, t) = \frac{p}{2t^{\frac{\zeta}{\alpha}p+1}} \langle\langle R, R \rangle\rangle + Cpt^{\left(\frac{\zeta}{\alpha}+1\right)p-1} f(X), \quad (3.10)$$

205 where $X \in \mathcal{Q}$ denotes position on \mathcal{Q} , V and R are velocity vector and momentum covector fields, t
 206 is the time variable, C is a constant, and ζ is given by equation (2.1).

207 **Theorem 3.3.** *The p -Bregman Euler–Lagrange equation corresponding to the p -Bregman La-*
 208 *grangian \mathcal{L}_p^{WQC} is given by*

$$209 \quad \nabla_{\dot{X}} \dot{X} + \frac{\zeta p + \alpha}{\alpha t} \dot{X} + Cp^2 t^{p-2} \text{grad}f(X) = 0. \quad (3.11)$$

210 *Proof.* See Appendix A.1, with $\lambda = \zeta/\alpha$.

211 **Theorem 3.4.** *Suppose $f : \mathcal{Q} \rightarrow \mathbb{R}$ is a geodesically α -weakly-quasi-convex function, and suppose*
 212 *that Assumption 1 is satisfied. Then, the p -Bregman Euler–Lagrange equation (3.11) has a solution,*
 213 *and any solution converges to a minimizer x^* of f with rate*

$$214 \quad f(X(t)) - f(x^*) \leq \frac{\zeta \|\text{Log}_{x_0}(x^*)\|^2}{2C\alpha^2 t^p}. \quad (3.12)$$

215 *Proof.* See Appendix B.1 for the existence of a solution and Appendix C.2 for the convergence rate.

216 **3.3. Strongly Convex Case.** Suppose $f : \mathcal{Q} \rightarrow \mathbb{R}$ is a geodesically μ -strongly-convex function,
 217 and suppose that Assumption 1 is satisfied. With ζ given by equation (2.1), let

$$218 \quad \eta = \left(\frac{1}{\sqrt{\zeta}} + \sqrt{\zeta} \right) \sqrt{\mu}. \quad (3.13)$$

219 We define the corresponding Bregman Lagrangian $\mathcal{L}^{SC} : T\mathcal{Q} \times \mathbb{R} \rightarrow \mathbb{R}$ via

$$220 \quad \mathcal{L}^{SC}(X, V, t) = \frac{e^{\eta t}}{2} \langle V, V \rangle - e^{\eta t} f(X), \quad (3.14)$$

221 and the corresponding Bregman Hamiltonian $\mathcal{H}^{SC} : T^*\mathcal{Q} \times \mathbb{R} \rightarrow \mathbb{R}$ is given by

$$222 \quad \mathcal{H}^{SC}(X, R, t) = \frac{e^{-\eta t}}{2} \langle\langle R, R \rangle\rangle + e^{\eta t} f(X). \quad (3.15)$$

223 **Theorem 3.5.** *The Bregman Euler–Lagrange equation corresponding to the Bregman Lagrangian*
 224 *\mathcal{L}^{SC} is given by*

$$225 \quad \nabla_{\dot{X}} \dot{X} + \eta \dot{X} + \text{grad}f(X) = 0. \quad (3.16)$$

226 *Proof.* The derivation of the Bregman Euler–Lagrange equation is presented in Appendix A.2.

227 **Theorem 3.6.** *Suppose $f : \mathcal{Q} \rightarrow \mathbb{R}$ is a geodesically μ -strongly-convex function, and suppose that*
228 *Assumption 1 is satisfied. Then, the Bregman Euler–Lagrange equation (3.16) has a solution, and*
229 *any solution converges to a minimizer x^* of f with rate*

$$230 \quad \boxed{f(X(t)) - f(x^*) \leq \frac{\mu \|\text{Log}_{x_0}(x^*)\|^2 + 2(f(x_0) - f(x^*))}{2e\sqrt{\frac{\mu}{\zeta}}t}}. \quad (3.17)$$

231 *Proof.* See Appendix B.2 for the existence and Theorem 7 from [4] for the convergence rate.

232

4. NUMERICAL EXPERIMENTS

233

234 The p -Bregman Euler–Lagrange equations can be rewritten as the first order system

$$235 \quad \dot{X} = V, \quad \nabla_V V = -\frac{\lambda p + 1}{t} V - Cp^2 t^{p-2} \text{grad}f(X), \quad (4.1)$$

236 where $\lambda = \zeta$ in the geodesically convex case and $\lambda = \zeta/\alpha$ in the geodesically α -weakly-quasi-convex
237 case, and as the first-order system

$$238 \quad \dot{X} = V, \quad \nabla_V V = -\left(\frac{1}{\sqrt{\zeta}} + \sqrt{\zeta}\right) \sqrt{\mu} V - \text{grad}f(X), \quad (4.2)$$

239 for the μ -strongly convex case. As in [4], we can adapt a semi-implicit Euler scheme (explicit Euler
240 update for the velocity V followed by an update for position X based on the updated value of V)
241 to the Riemannian setting to obtain the following algorithm:

Algorithm 1: Semi-Implicit Euler Integration of the p -Bregman Euler–Lagrange Equations

Input: A function $f : \mathcal{Q} \rightarrow \mathbb{R}$. Constants $C, h, p > 0$. $X_0 \in \mathcal{Q}$. $V_0 \in T_{X_0} \mathcal{Q}$.

```

1 while convergence criterion is not met do
2   if  $f$  is  $\mu$ -geodesically strongly convex then
3      $b_k \leftarrow 1 - h\left(\frac{1}{\sqrt{\zeta}} + \sqrt{\zeta}\right)\sqrt{\mu}$ ,  $c_k \leftarrow 1$ 
242 4   else if  $f$  is geodesically convex ( $\lambda = \zeta$ ) or  $\alpha$ -weakly-quasi-convex ( $\lambda = \zeta/\alpha$ ) then
5      $b_k \leftarrow 1 - \frac{\lambda p + 1}{k}$ ,  $c_k \leftarrow Cp^2(kh)^{p-2}$ 
6     Version I:  $a_k \leftarrow b_k V_k - hc_k \text{grad}f(X_k)$ 
7     Version II:  $a_k \leftarrow b_k V_k - hc_k \text{grad}f(\text{Exp}_{X_k}(hb_k V_k))$ 
8    $X_{k+1} \leftarrow \text{Exp}_{X_k}(ha_k)$ ,  $V_{k+1} \leftarrow \Gamma_{X_k}^{X_{k+1}} a_k$ 

```

243 Version I of Algorithm 1 corresponds to the usual update for the Semi-Implicit Euler scheme,
244 while Version II is inspired by the reformulation of Nesterov’s method from [23] that uses a cor-
245 rected gradient $\nabla f(X_k + hb_k V_k)$ instead of the traditional gradient $\nabla f(X_k)$. Note that the SIRMAG
246 algorithm presented in [4] corresponds to the special case where $p = 2$ and $C = 1/4$.

247

248 The first problem we have investigated is the problem presented in [4] of minimizing the (strongly
249 convex) distance function $f(x) = \frac{1}{2}d(x, q)^2$ for a given point q , on a subset of chosen finite diameter
250 of the hyperbolic plane \mathbb{H}^2 , which is a manifold with constant negative curvature $K = -1$.

251 The second problem we have investigated is Rayleigh quotient optimization. Eigenvectors corre-
252 sponding to the largest eigenvalue of a symmetric $n \times n$ matrix A maximize the Rayleigh quotient
253 $\frac{v^T A v}{v^T v}$ over \mathbb{R}^n . Thus, a unit eigenvector v^* corresponding to the largest eigenvalue of the matrix A
254 is a minimizer of the function $f(v) = -v^T A v$, over the unit sphere $\mathcal{Q} = \mathbb{S}^{n-1}$, which can be thought
255 of as a Riemannian submanifold with constant positive curvature $K = 1$ of \mathbb{R}^n endowed with the
256 Riemannian metric inherited from the Euclidean inner product $g_v(u, w) = u^T w$. More information

257 concerning the geometry of \mathbb{S}^{n-1} , such as its tangent bundle, its orthogonal projection and expo-
 258 nential map can be found in [1]. Solving the Rayleigh quotient optimization problem efficiently
 259 is challenging when the given symmetric matrix A is ill-conditioned and high-dimensional. Note
 260 that an efficient algorithm that solves the above minimization problem can also be used to find
 261 eigenvectors corresponding to the smallest eigenvalue of A by using the fact that the eigenvalues of
 262 A are the negative of the eigenvalues of $-A$.

263

264 Experiments carried out in [4] showed that SIRNAG (the convex $p = 2$ Algorithm 1) and the
 265 strongly convex Algorithm 1 were of comparable efficiency or more efficient than the standard Rie-
 266 mannian Gradient Descent (RGD) method, depending on the properties of the objective function
 267 and on the geometry of the Riemannian manifold. We have conducted further numerical experi-
 268 ments to investigate how the simple discretization of higher-order $p = 6$ Bregman dynamics com-
 269 pared to its $p = 2$ counterpart, and to see whether it matches the theoretical $\mathcal{O}(t^{-p})$ conver-
 270 gence rate. The numerical results obtained for the distance minimization and Rayleigh minimization
 271 problems are illustrated in Figure 1, where all the algorithms were implemented with the same
 272 fixed timestep. We can see that the $p = 6$ algorithms outperform their $p = 2$ counterparts, and that
 273 the efficiency improvement is very important. Furthermore, both versions of the $p = 6$ Algorithm 1
 274 exhibit a faster convergence rate than the theoretical $\mathcal{O}(t^{-6})$ rate. While Version I of Algorithm 1
 275 exhibits polynomial rates of $\mathcal{O}(t^{-10.8})$ and $\mathcal{O}(t^{-9})$ on the objective functions considered, Version II
 276 of Algorithm 1 exhibits a much faster exponential rate of convergence on both examples.

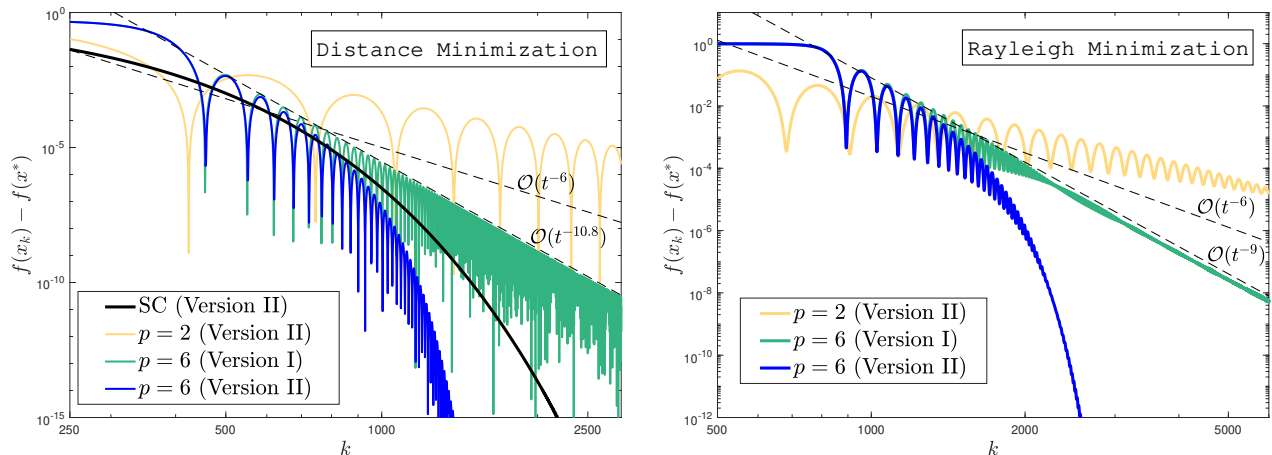


FIGURE 1. Comparison of the rates of convergence of the μ -strongly convex (SC) Algorithm 1 and convex Algorithms 1 with different values of p and with the two versions of the update corresponding to the traditional and corrected gradients. Note that all the algorithms were implemented with the same timestep h .

277 Note however that an increase in the value of p in Algorithm 1, which corresponds to an increase
 278 in the order of the Bregman dynamics integrated, requires a decrease in the timestep, in agreement
 279 with intuitive expectations. This timestep decrease requirement is especially important due to
 280 the polynomially growing $h(kh)^{p-2}$ coefficient multiplying the gradient of f in the updates of the
 281 algorithm. Such a decrease in the timestep does not really affect the convergence rate, but the
 282 transition between the initialization and convergence phases takes longer.

283 Similar issues arise when discretizing the continuous Euler–Lagrange flow associated with ac-
 284 celerated optimization on vector spaces, and in that situation, it was observed that time-adaptive

285 symplectic integrators based on Hamiltonian variational integrators resulted in dramatically im-
286 proved robustness and stability. As such, it will be natural to explore generalizations of time-
287 adaptive symplectic integrators based on Hamiltonian variational integrators applied to Poincaré
288 transformed Hamiltonians, that respect the Riemannian manifold structure in order to yield more
289 robust and stable numerical discretizations of the flows we have studied in this paper in order to
290 construct accelerated optimization algorithms on Riemannian manifolds.

291 Finally, Figure 2 shows that the discretization empirically converges to the solution of the ODE
292 as the timestep h goes to 0. Note that although all the discretizations follow the ODE trajectory
293 closely, smaller timesteps result in a larger number of iterations, especially to transition from the
294 initialization plateau to the convergence phase (around time $t = 4$ in the example presented in
295 Figure 2). A theoretical shadowing result bounding the error between the discrete-time RGD and
296 its continuous-time limiting ODE was obtained in [4]. It would be desirable to obtain similar
297 shadowing results in the future for discretizations of the class of ODEs considered here, perhaps
298 drawing inspiration from [27]. However, such a result might be very difficult to obtain because
299 momentum methods lack contraction, are nondescending, and are highly oscillatory [4; 21].

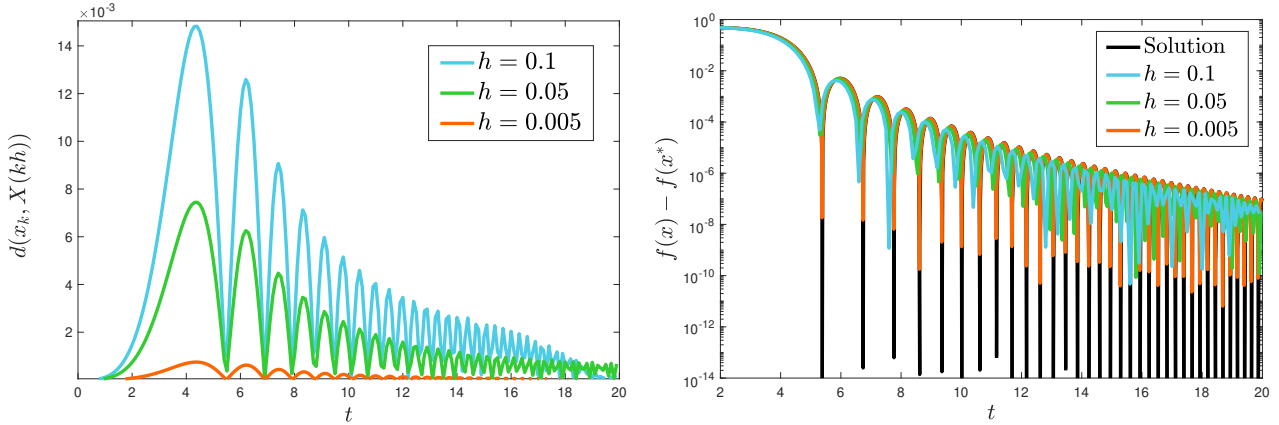


FIGURE 2. Discretization errors (top graph) and convergence rates (bottom graphs) of Version I of the $p = 5$ convex Algorithm 1 with different values of h for the distance minimization problem. The true solution of the differential equation was approximated by the same algorithm with a very small timestep $h = 10^{-5}$.

5. TIME INVARIANCE AND POINCARÉ TRANSFORMATION

300 Let $f : \mathcal{Q} \rightarrow \mathbb{R}$ be a given function, and suppose that Assumption 1 is satisfied. In both the
302 cases where f is geodesically convex and α -weakly-quasi-convex, we have formulated in section 3
303 a variational framework for the minimization of f , via a p -Bregman Lagrangian $\mathcal{L}_p : T\mathcal{Q} \times \mathbb{R} \rightarrow \mathbb{R}$
304 and a corresponding p -Bregman Hamiltonian $\mathcal{H}_p : T^*\mathcal{Q} \times \mathbb{R} \rightarrow \mathbb{R}$ for $p > 0$ of the form

$$\mathcal{L}_p(X, V, t) = \frac{t^{\lambda p + 1}}{2p} \langle V, V \rangle - C p t^{(\lambda + 1)p - 1} f(X), \quad (5.1)$$

$$\mathcal{H}_p(X, R, t) = \frac{p}{2t^{\lambda p + 1}} \langle R, R \rangle + C p t^{(\lambda + 1)p - 1} f(X), \quad (5.2)$$

308 with associated p -Bregman Euler–Lagrange equations given by

$$\nabla_{\dot{X}} \dot{X} + \frac{\lambda p + 1}{t} \dot{X} + C p^2 t^{p-2} \text{grad}f(X) = 0, \quad (5.3)$$

where $\lambda = \zeta$ in the geodesically convex case, and $\lambda = \zeta/\alpha$ in the geodesically α -weakly-quasi-convex case. Theorems 3.2 and 3.4 imply that in both cases, solutions to the p -Bregman Euler–Lagrange equations converge to a minimizer of f at a convergence rate of $\mathcal{O}(1/t^p)$. Now, the following two theorems show that in both cases, time-rescaling via $\tau(t) = t^{\mathring{p}/p}$ a solution to the p -Bregman Euler–Lagrange equations yields a solution to the \mathring{p} -Bregman Euler–Lagrange equations.

Theorem 5.1. *Suppose $f : \mathcal{Q} \rightarrow \mathbb{R}$ is a geodesically convex function, and Assumption 1 is satisfied. Suppose the curve $X(t)$ satisfies the corresponding p -Bregman Euler–Lagrange equation (3.7). Then, the reparametrized curve $X(t^{\mathring{p}/p})$ satisfies the \mathring{p} -Bregman Euler–Lagrange equation (3.7).*

Proof. See Appendix D with $\lambda = \zeta$.

Theorem 5.2. *Suppose $f : \mathcal{Q} \rightarrow \mathbb{R}$ is geodesically α -weakly-quasi-convex function, and Assumption 1 is satisfied. Suppose the curve $X(t)$ satisfies the p -Bregman Euler–Lagrange equation (3.11). Then, the reparametrized curve $X(t^{\mathring{p}/p})$ satisfies the \mathring{p} -Bregman Euler–Lagrange equation (3.11).*

Proof. See Appendix D with $\lambda = \zeta/\alpha$.

Thus, the entire subfamily of Bregman trajectories indexed by the parameter p can be obtained by speeding up or slowing down along the Bregman curve in spacetime corresponding to any specific value of p . Inspired by the computational efficiency of the approach introduced in [7], it is natural to attempt to exploit the time-rescaling property of the Bregman dynamics together with a carefully chosen Poincaré transformation to transform the p -Bregman Hamiltonian into an autonomous version of the \mathring{p} -Bregman Hamiltonian in extended phase-space, where $\mathring{p} < p$. This would allow us to integrate the higher-order p -Bregman dynamics while benefiting from the computational efficiency of integrating the lower-order \mathring{p} -Bregman dynamics. Explicitly, the time rescaling $\tau(t) = t^{\mathring{p}/p}$ is associated to the monitor function

$$\frac{dt}{d\tau} = g_{p \rightarrow \mathring{p}}(t) = \frac{p}{\mathring{p}} t^{1-\mathring{p}/p}, \quad (5.4)$$

and generates a Poincaré transformed Hamiltonian

$$\bar{\mathcal{H}}_{p \rightarrow \mathring{p}}(\bar{X}, \bar{R}) = g_{p \rightarrow \mathring{p}}(X^t) (\mathcal{H}_p(\bar{X}, R) + R^t), \quad (5.5)$$

in the extended space $\bar{\mathcal{Q}} = \mathcal{Q} \times \mathbb{R}$ where $\bar{X} = \begin{bmatrix} X \\ X^t \end{bmatrix}$ and $\bar{R} = \begin{bmatrix} R \\ R^t \end{bmatrix}$. We will make the conventional choice $X^t = t$ and R^t as the conjugate momentum of X^t with $R^t(0) = -\mathcal{H}_p(X(0), R(0), 0) = -H_0$, chosen so that $\bar{\mathcal{H}}_{p \rightarrow \mathring{p}}(\bar{X}, \bar{R}) = 0$ along all integral curves through $(\bar{X}(0), \bar{R}(0))$. The time t shall be referred to as the physical time, while τ will be referred to as the fictive time. The corresponding Hamiltonian equations of motion in the extended phase space are then given by

$$\dot{\bar{X}} = \frac{\partial \bar{\mathcal{H}}_{p \rightarrow \mathring{p}}}{\partial \bar{R}}, \quad \dot{\bar{R}} = -\frac{\partial \bar{\mathcal{H}}_{p \rightarrow \mathring{p}}}{\partial \bar{X}}. \quad (5.6)$$

Now, suppose $(\bar{X}(\tau), \bar{R}(\tau))$ are solutions to these extended equations of motion, and let $(x(t), r(t))$ solve Hamilton’s equations for the original Hamiltonian \mathcal{H}_p . Then

$$\bar{\mathcal{H}}_{p \rightarrow \mathring{p}}(\bar{X}(\tau), \bar{R}(\tau)) = \bar{\mathcal{H}}_{p \rightarrow \mathring{p}}(\bar{X}(0), \bar{R}(0)) = 0.$$

Thus, the components $(X(\tau), R(\tau))$ in the original phase space of $(\bar{X}(\tau), \bar{R}(\tau))$ satisfy

$$\mathcal{H}_p(X(\tau), R(\tau), \tau) = -R^t(\tau), \quad \mathcal{H}_p(X(0), R(0), 0) = -R^t(0) = \mathcal{H}_p(x(0), r(0), 0).$$

Therefore, $(X(\tau), R(\tau))$ and $(x(t), r(t))$ both satisfy Hamilton’s equations for the original Hamiltonian \mathcal{H}_p with the same initial values, so they must be the same.

346 As a consequence, instead of integrating the p -Bregman Hamiltonian system (5.2), we can focus
 347 on the Poincaré transformed Hamiltonian $\bar{\mathcal{H}}_{p \rightarrow \dot{p}}$ in extended phase-space given by equation (5.5),
 348 with \mathcal{H}_p and $g_{p \rightarrow \dot{p}}$ given by equations (5.2) and (5.4), that is

$$349 \quad \bar{\mathcal{H}}_{p \rightarrow \dot{p}}(\bar{X}, \bar{R}) = \frac{p^2}{2\dot{p}(X^t)^{\lambda+\dot{p}/p}} \langle\langle R, R \rangle\rangle + \frac{Cp^2}{\dot{p}} (X^t)^{(\lambda+1)p-\dot{p}/p} f(X) + \frac{p}{\dot{p}} (X^t)^{1-\dot{p}/p} R^t, \quad (5.7)$$

350 where $\lambda = \zeta$ if f is geodesically convex, and $\lambda = \zeta/\alpha$ if f is geodesically α -weakly-quasi-convex. The
 351 resulting integrator has constant timestep in fictive time τ but variable timestep in physical time t .

352 In our prior work on discretizations of variational formulations of accelerated optimization on
 353 normed spaces [7], we performed a very careful computational study of how time-adaptivity and
 354 symplecticity of the numerical scheme improve the performance of the resulting numerical optimiza-
 355 tion algorithm. In particular, we observed that time-adaptive Hamiltonian variational discretiza-
 356 tions, which are automatically symplectic, with adaptive timesteps informed by the time invariance
 357 of the family of p -Bregman Lagrangians and Hamiltonians yielded the most robust and computa-
 358 tionally efficient numerical optimization algorithms, outperforming fixed-timestep symplectic dis-
 359 cretizations, adaptive-timestep non-symplectic discretizations, and Nesterov’s accelerated gradient
 360 algorithm which is neither time-adaptive nor symplectic. As such, it would be desirable to general-
 361 ize the time-adaptive Hamiltonian variational integrator framework to Riemannian manifolds, and
 362 apply it to the variational formulation of accelerated optimization on Riemannian manifolds.

363

364

6. CONCLUSION

365 We have shown that on Riemannian manifolds, the convergence rate in continuous time of a
 366 geodesically convex, α -weakly-quasi convex, or μ -strongly convex function $f(x(t))$ to its opti-
 367 mal value can be accelerated to an arbitrary convergence rate $\mathcal{O}(1/t^p)$, which extended the re-
 368 sults of [24] from normed vector spaces to Riemannian manifolds. This rate of convergence is
 369 achieved along solutions of the Euler–Lagrange and Hamilton’s equations corresponding to a fam-
 370 ily of time-dependent Bregman Lagrangian and Hamiltonian systems on Riemannian manifolds. As
 371 was demonstrated in the normed vector space setting, such families of Bregman Lagrangians and
 372 Hamiltonians can be used to construct practical, robust, and computationally efficient numerical
 373 optimization algorithms that outperform Nesterov’s accelerated gradient method by considering
 374 geometric structure-preserving discretizations of the continuous-time flows.

375 Numerical experiments implementing a simple discretization of the p -Bregman Euler–Lagrange
 376 equations applied to a distance minimization and Rayleigh minimization problems confirmed that
 377 the higher-order algorithms outperform significantly their lower-order counterparts and their the-
 378 oretical $\mathcal{O}(t^{-p})$ convergence rates. Numerical results also showed that using a corrected gradient
 379 in the update instead of the traditional gradient, as was done in [23], improved the theoretically
 380 predicted polynomial convergence rate to an exponential rate of convergence in practice. While
 381 higher values of p result in faster rates of convergence, they also appear to be more prone to sta-
 382 bility issues under numerical discretization, which can cause the numerical optimization algorithm
 383 to diverge, but we anticipate that symplectic discretizations will address these stability issues.

384 Finally, in analogy to what was done in [24] for normed vector spaces, we proved that the family
 385 of time-dependent Bregman Lagrangian and Hamiltonians on Riemannian manifolds is closed under
 386 time rescaling. Inspired by the computational efficiency of the approach introduced in [7], we can
 387 then exploit this invariance property via a carefully chosen Poincaré transformation that will allow
 388 us to integrate higher-order p -Bregman dynamics while benefiting from the computational efficiency
 389 of integrating a lower-order \dot{p} -Bregman Hamiltonian system.

390 It was observed in our prior computational experiments in the normed vector space case [7]
 391 that geometric discretizations which respect the time-rescaling invariance and symplecticity of the

392 Bregman Lagrangian and Hamiltonian flows were substantially less prone to stability issues, and
 393 were therefore more robust, reliable, and computationally efficient. As such, it is natural to develop
 394 time-adaptive Hamiltonian variational integrators for the Bregman Hamiltonian introduced in this
 395 paper describing accelerated optimization on Riemannian manifolds.

396 Developing an intrinsic extension of Hamiltonian variational integrators to manifolds will require
 397 some additional work, since the current approach involves Type II/Type III generating functions
 398 $H_d^+(q_k, p_{k+1})$, $H_d^-(p_k, q_{k+1})$, which depend on the position at one boundary point, and the momen-
 399 tum at the other boundary point. However, this does not make intrinsic sense on a manifold, since
 400 one needs the base point in order to specify the corresponding cotangent space, and one should
 401 ideally consider a Hamiltonian variational integrator construction based on discrete Dirac mechan-
 402 ics [13], which would yield a generating function $E_d^+(q_k, q_{k+1}, p_{k+1})$, $E_d^-(q_k, p_k, q_{k+1})$, that depends
 403 on the position at both boundary points and the momentum at one of the boundary points. This
 404 approach can be viewed as a discretization of the generalized energy $E(q, v, p) = \langle p, v \rangle - L(q, v)$, in
 405 contrast to the Hamiltonian $H(q, p) = \text{ext}_v \langle p, v \rangle - L(q, v) = \langle p, v \rangle - L(q, v)|_{p=\frac{\partial L}{\partial v}}$.

406 However, a more practical method relies on the fact that we have a Riemannian manifold, which
 407 is endowed with a Riemannian exponential and Riemannian logarithm that can be used to construct
 408 an extension of Hamiltonian variational integrators using geodesic normal coordinates. For many
 409 important matrix manifolds, one can replace the Riemannian exponential in the geodesic normal
 410 coordinates by a retraction [1], which is often constructed using matrix factorizations.

411 Another important case involves Riemannian submanifolds that are embedded in a Riemannian
 412 linear manifold and are realized as the level set of a submersion. The characterization of the sub-
 413 manifold as the level set of a submersion, together with the linear space structure of the embedding
 414 space, and the variational characterization of the dynamics naturally lends itself to the use of the
 415 Lagrange multiplier theorem, which allows one to use Hamiltonian variational integrators defined
 416 on the embedding space by including a Lagrange multiplier term involving the submersion in the
 417 Lagrangian or Hamiltonian. This is analogous to the derivation of the SHAKE and RATTLE
 418 methods as a variational integrator for constrained systems (see, for example, §3.5 of [16]).

419 We anticipate that applying an appropriate generalization of Hamiltonian variational integrators
 420 to the Bregman Hamiltonians introduced in this paper will yield a novel class of robust and efficient
 421 accelerated optimization algorithms on Riemannian manifolds, and we intend to pursue this research
 422 direction in future work.

423 ACKNOWLEDGMENTS

424 The authors were supported in part by NSF under grants DMS-1411792, DMS-1345013, DMS-
 425 1813635, and by AFOSR under grant FA9550-18-1-0288.

427 APPENDIX A. DERIVATION OF THE EULER–LAGRANGE EQUATIONS

428 A.1. Convex and Weakly-Quasi-Convex Cases.

429 **Theorem A.1.** *The Euler–Lagrange equation corresponding to the Lagrangian*

$$430 \quad \mathcal{L}(X, V, t) = \frac{t^{\lambda p + 1}}{2p} \langle V, V \rangle - C p t^{(\lambda + 1)p - 1} f(X),$$

431 *is given by*

$$432 \quad \nabla_{\dot{X}} \dot{X} + \frac{\lambda p + 1}{t} \dot{X} + C p^2 t^{p-2} \text{grad} f(X) = 0.$$

433 *Proof.* Consider a path on the manifold \mathcal{Q} described in coordinates by

$$434 \quad (x(t), \dot{x}(t)) = (q^1(t), \dots, q^n(t), v^1(t), \dots, v^n(t)).$$

435 Then, with $\langle \cdot, \cdot \rangle = \sum_{i,j=1}^n g_{ij} dx^i dx^j$, the p -Bregman Lagrangian can be written as

$$436 \quad \mathcal{L}(x(t), \dot{x}(t), t) = \frac{t^{\lambda p+1}}{2p} \sum_{i,j=1}^n g_{ij}(x(t)) v^i(t) v^j(t) - C p t^{(\lambda+1)p-1} f(x(t)).$$

437 For $k = 1, \dots, n$,

$$438 \quad \frac{d}{dt} \left(\frac{\partial \mathcal{L}}{\partial v^k} (x(t), \dot{x}(t), t) \right) = \frac{t^{\lambda p+1}}{p} \sum_{i=1}^n g_{ik}(x(t)) \frac{dv^i}{dt}(t) + \frac{t^{\lambda p+1}}{p} \sum_{i,j=1}^n \frac{\partial g_{kj}}{\partial q^i} (x(t)) v^i(t) v^j(t) \\ 439 \quad \quad \quad + \frac{\lambda p+1}{p} t^{\lambda p} \sum_{i=1}^n g_{ik}(x(t)) v^i(t),$$

$$440 \quad \frac{\partial \mathcal{L}}{\partial q^k} (x(t), \dot{x}(t), t) = \frac{t^{\lambda p+1}}{2p} \sum_{i,j=1}^n \frac{\partial g_{ij}}{\partial q^k} (x(t)) v^i(t) v^j(t) - C p t^{(\lambda+1)p-1} \frac{\partial f}{\partial q^k} (x(t)).$$

442 If we multiply both terms by $\frac{p}{t^{\lambda p+1}}$, the Euler–Lagrange equations (2.3) are given for $k = 1, \dots, n$ by

$$443 \quad 0 = \sum_{i=1}^n g_{ik}(x(t)) \frac{dv^i}{dt}(t) + \sum_{i,j=1}^n \frac{\partial g_{kj}}{\partial q^i} (x(t)) v^i(t) v^j(t) + \frac{\lambda p+1}{t} \sum_{i=1}^n g_{ik}(x(t)) v^i(t) \\ 444 \quad \quad \quad - \frac{1}{2} \sum_{i,j=1}^n \frac{\partial g_{ij}}{\partial q^k} (x(t)) v^i(t) v^j(t) + C p^2 t^{p-2} \frac{\partial f}{\partial q^k} (x(t)).$$

445 Multiplying by the matrix (g^{ij}) , which is the inverse of (g_{ij}) , we get for $k = 1, \dots, n$

$$446 \quad \left(\frac{dv^k}{dt}(t) + \sum_{i,j=1}^n \Gamma_{ij}^k(x(t)) v^i(t) v^j(t) \right) + \frac{\lambda p+1}{t} v^k(t) + C p^2 t^{p-2} (\text{grad} f(x(t)))^k = 0,$$

447 where Γ_{ij}^k are the Christoffel symbols given by $\Gamma_{ij}^k = \frac{1}{2} \sum_{l=1}^n g^{kl} \left[\frac{\partial g_{jl}}{\partial x^i} + \frac{\partial g_{li}}{\partial x^j} - \frac{\partial g_{ij}}{\partial x^l} \right]$, which gives the
448 desired Euler–Lagrange equation. \square

449

450 A.2. Strongly Convex Case.

451 **Theorem A.2.** *The Bregman Euler–Lagrange equation corresponding to the Bregman Lagrangian*
452 \mathcal{L}^{SC} *is given by*

$$453 \quad \boxed{\nabla_{\dot{X}} \dot{X} + \eta \dot{X} + \text{grad} f(X) = 0.}$$

454 *Proof.* Consider a path on the manifold \mathcal{Q} described in coordinates by

$$455 \quad (x(t), \dot{x}(t)) = (q^1(t), \dots, q^n(t), v^1(t), \dots, v^n(t)).$$

456 Then, with $\langle \cdot, \cdot \rangle = \sum_{i,j=1}^n g_{ij} dx^i dx^j$, the Bregman Lagrangian can be written as

$$457 \quad \mathcal{L}_p^{SC}(x(t), \dot{x}(t), t) = \frac{e^{\eta t}}{2} \sum_{i,j=1}^n g_{ij}(x(t)) v^i(t) v^j(t) - e^{\eta t} f(x(t)).$$

458 For $k = 1, \dots, n$

$$459 \quad \frac{d}{dt} \left(\frac{\partial \mathcal{L}_p^{SC}}{\partial v^k} (x(t), \dot{x}(t), t) \right) = e^{\eta t} \sum_{i=1}^n g_{ik}(x(t)) \frac{dv^i}{dt}(t) + e^{\eta t} \sum_{i,j=1}^n \frac{\partial g_{kj}}{\partial q^i} (x(t)) v^i(t) v^j(t) \\ 460 \quad \quad \quad + \eta e^{\eta t} \sum_{i=1}^n g_{ik}(x(t)) v^i(t),$$

$$461 \quad \frac{\partial \mathcal{L}_p^{SC}}{\partial q^k} (x(t), \dot{x}(t), t) = e^{\eta t} \sum_{i,j=1}^n \frac{\partial g_{ij}}{\partial q^k} (x(t)) v^i(t) v^j(t) - e^{\eta t} \frac{\partial f}{\partial q^k} (x(t)).$$

462

463 If we multiply both terms by $e^{-\eta t}$, the Euler–Lagrange equations (2.3) are given for $k = 1, \dots, n$ by

$$464 \quad 0 = \sum_{i=1}^n g_{ik}(x(t)) \frac{dv^i}{dt}(t) + \sum_{i,j=1}^n \frac{\partial g_{kj}}{\partial q^i}(x(t)) v^i(t) v^j(t) + \eta \sum_{i=1}^n g_{ik}(x(t)) v^i(t) \\ 465 \quad - \frac{1}{2} \sum_{i,j=1}^n \frac{\partial g_{ij}}{\partial q^k}(x(t)) v^i(t) v^j(t) + \frac{\partial f}{\partial q^k}(x(t)).$$

466 Rearranging terms, and multiplying by the matrix (g^{ij}) which is the inverse of (g_{ij}) , we get for
467 $k = 1, \dots, n$

$$468 \quad \left(\frac{dv^k}{dt}(t) + \sum_{i,j=1}^n \Gamma_{ij}^k(x(t)) v^i(t) v^j(t) \right) + \eta v^k(t) + (\text{grad}f(x(t)))^k = 0,$$

469 where Γ_{ij}^k are the Christoffel symbols given by $\Gamma_{ij}^k = \frac{1}{2} \sum_{l=1}^n g^{kl} \left[\frac{\partial g_{jl}}{\partial x^i} + \frac{\partial g_{li}}{\partial x^j} - \frac{\partial g_{ij}}{\partial x^l} \right]$, which gives the
470 desired Euler–Lagrange equation. \square

471

472

APPENDIX B. PROOF OF EXISTENCE THEOREMS

473 B.1. Convex and α -Weakly-Quasi-Convex Cases.

474 **Theorem B.1.** *Suppose Assumption 1 is satisfied, and let $C, p > 0$ and $v > 1$ be given constants.*
475 *Then the differential equation*

$$476 \quad \nabla_{\dot{X}} \dot{X} + \frac{v}{t} \dot{X} + Ct^{p-2} \text{grad}f(X) = 0,$$

477 *has a global solution $X : [0, \infty) \rightarrow \mathcal{Q}$ under the initial conditions $X(0) = x_0 \in \mathcal{Q}$ and $\dot{X}(0) = 0$.*

478 *Proof.* The proof is similar to that of Lemma 3 in [4], which extended Theorem 1 in [22] to the
479 Riemannian setting. We first define a family of smoothed equation for which we then show existence
480 of a solution for all time. After choosing an equicontinuous and uniformly bounded subfamily of
481 smoothed solutions, we use the Arzela–Ascoli Theorem on the complete Riemannian manifold \mathcal{Q} to
482 obtain a subsequence converging uniformly, and argue that the limit of this subsequence is solves
483 the original problem. When $p = 2$, we recover the simpler case considered in Lemma 3 of [4], so we
484 assume $p \neq 2$ in this proof. Consider the following families of smoothed equations for $\delta > 0$:

$$485 \quad \nabla_{\dot{X}} \dot{X} + \frac{v}{\max(\delta, t)} \dot{X} + C(\max(\delta, t))^{p-2} \text{grad}f(X) = 0 \quad \text{if } p < 2,$$

$$486 \quad \nabla_{\dot{X}} \dot{X} + \frac{v}{\max(\delta, t)} \dot{X} + Ct^{p-2} \text{grad}f(X) = 0 \quad \text{if } p > 2.$$

487 Exp and Log are defined globally on \mathcal{Q} by Assumption 1, so we can choose geodesically normal
488 coordinates $\phi = \psi^{-1}$ around x_0 defined globally on \mathcal{Q} and put $c = \phi \circ X$. Using the smoothness of
489 f and letting $u = \dot{c}$ gives a system of first order ODEs defining a local representation for a vector
490 field in $T\mathcal{Q}$, and section IV.3 of [11] guarantees that the smoothed ODE has a unique solution X_δ
491 locally around 0. Actually, X_δ exists on $[0, \infty)$. Indeed, by contradiction, let $[0, T)$ be the maximal
492 interval of existence of X_δ , for some finite $T > 0$. Using $\frac{d}{dt} f(X_\delta(t)) = \langle \text{grad}f(X_\delta), \dot{X}_\delta \rangle$ gives

$$493 \quad \frac{d}{dt} f(X_\delta) = -\frac{\delta^{2-p}}{C} \langle \nabla_{\dot{X}_\delta} \dot{X}_\delta, \dot{X}_\delta \rangle - \frac{v\delta^{1-p}}{C} \langle \dot{X}_\delta, \dot{X}_\delta \rangle = -\frac{\delta^{2-p}}{2C} \frac{d}{dt} \|\dot{X}_\delta\|^2 - \frac{v\delta^{1-p}}{C} \|\dot{X}_\delta\|^2 \quad \text{if } \delta > t, p < 2,$$

$$494 \quad \frac{d}{dt} f(X_\delta) = -\frac{t^{2-p}}{C} \langle \nabla_{\dot{X}_\delta} \dot{X}_\delta, \dot{X}_\delta \rangle - \frac{vt^{2-p}}{C\delta} \langle \dot{X}_\delta, \dot{X}_\delta \rangle = -\frac{t^{2-p}}{2C} \frac{d}{dt} \|\dot{X}_\delta\|^2 - \frac{vt^{2-p}}{C\delta} \|\dot{X}_\delta\|^2 \quad \text{if } \delta > t, p > 2,$$

$$495 \quad \frac{d}{dt} f(X_\delta) = -\frac{t^{2-p}}{C} \langle \nabla_{\dot{X}_\delta} \dot{X}_\delta, \dot{X}_\delta \rangle - \frac{vt^{1-p}}{C} \langle \dot{X}_\delta, \dot{X}_\delta \rangle = -\frac{1}{2C} \frac{d}{dt} (t^{2-p} \|\dot{X}_\delta\|^2) - \frac{2v(2-p)-1}{2C(2-p)} t^{1-p} \|\dot{X}_\delta\|^2 \quad \text{if } \delta < t.$$

496 Let $\theta = \frac{2v(2-p)-1}{2C(2-p)}$. Integrating and using the Cauchy-Schwarz inequality for the $p < 2$ case gives

$$\begin{aligned}
 497 \quad & \int_0^T \sqrt{(\max(\delta, t))^{1-p}} \|\dot{X}_\delta\| dt = \int_0^\delta \sqrt{\delta^{1-p}} \|\dot{X}_\delta\| dt + \int_\delta^T \sqrt{t^{1-p}} \|\dot{X}_\delta\| dt \\
 498 \quad & \leq \sqrt{\frac{C\delta}{v} (f(x_0) - \inf_u f(u)) + \frac{\delta^{2-p}}{2v} \left(\|\dot{X}_\delta(0)\|^2 - \inf_{t \in [0, T]} \|\dot{X}_\delta(t)\|^2 \right)} \\
 499 \quad & \quad + \sqrt{\frac{T-\delta}{\theta} (f(X_\delta(\delta)) - \inf_u f(u)) + \frac{T-\delta}{2C\theta} \left(\delta^{2-p} \|\dot{X}_\delta(\delta)\|^2 - \inf_{t \in [0, T]} t^{2-p} \|\dot{X}_\delta(t)\|^2 \right)} < \infty,
 \end{aligned}$$

500 since f is bounded below by Assumption 1. If $\delta \geq T$, then $\sqrt{\delta^{1-p}} \dot{X}_\delta$ is integrable on $[0, T)$. If
 501 $\delta < T$, then the integrals on $[0, T)$ and $[0, \delta)$ are finite, so the integral on $[\delta, T)$ must also be fi-
 502 nite, so $\sqrt{t^{1-p}} \dot{X}_\delta$ is integrable on $[\delta, T)$. Now, $\|\int_a^T \dot{X}_\delta dt\| \leq \int_a^T \|\dot{X}_\delta\| dt < \infty$ for $a = 0, \delta$ implies that
 503 $\lim_{t \rightarrow T} X_\delta(t)$ exists. Since \mathcal{Q} is complete by Assumption 1, the limit is in \mathcal{Q} , contradicting the max-
 504 imality of $[0, T)$. The $p > 2$ case is similar: the integrand is replaced by $\sqrt{t^{2-p}(\max(\delta, t))^{-1}} \|\dot{X}_\delta\|$,
 505 and the integral on $[\delta, T)$ remains unchanged while the integral on $[0, \delta)$ can be bounded by the
 506 same expression using $t < \delta$. Thus, in both cases, we can find a solution $X_\delta : [0, \infty) \rightarrow \mathcal{Q}$ to the
 507 smooth initial-valued ODE, and its corresponding solution $X_\delta : [0, \infty) \rightarrow \mathbb{R}^n$ in local coordinates.
 508 Now let

$$509 \quad M_\delta(t) = \sup_{u \in (0, t]} \frac{\|\dot{X}_\delta(u)\|}{u}$$

510 When $0 < t \leq \delta$, the smoothed ODE can be written as

$$511 \quad \nabla_{\dot{X}_\delta} \left(\dot{X}_\delta e^{\frac{v}{\delta}} \right) = -C\delta^{p-2} \text{gradf}(X_\delta) e^{\frac{v}{\delta}} \text{ if } p < 2, \quad \nabla_{\dot{X}_\delta} \left(\dot{X}_\delta e^{\frac{v}{\delta}} \right) = -Ct^{p-2} \text{gradf}(X_\delta) e^{\frac{v}{\delta}} \text{ if } p > 2.$$

512 Thus, we can use Lemma 4 in [4] to get for $p > 2$ that

$$\begin{aligned}
 513 \quad & \Gamma_{X_\delta(t)}^{x_0} \dot{X}_\delta(t) = -e^{-\frac{v}{\delta}t} \int_0^t \left(\Gamma_{X_\delta(u)}^{x_0} \text{gradf}(X_\delta(u)) - \Gamma_{X_\delta(u)}^{x_0} \Gamma(X_\delta)_{x_0}^{X_\delta(u)} \text{gradf}(x_0) \right) C u^{p-2} e^{\frac{v}{\delta}u} du \\
 514 \quad & \quad - e^{-\frac{v}{\delta}t} \int_0^t C u^{p-2} \Gamma_{X_\delta(u)}^{x_0} \Gamma(X_\delta)_{x_0}^{X_\delta(u)} \text{gradf}(x_0) e^{\frac{v}{\delta}u} du.
 \end{aligned}$$

515 From the Lipschitz assumption on f , we have that

$$516 \quad \|\text{gradf}(X_\delta(u)) - \Gamma_{x_0}^{X_\delta(u)} \text{gradf}(x_0)\| \leq L \int_0^u \|\dot{X}_\delta(s)\| ds = L \int_0^u s \frac{\|\dot{X}_\delta(s)\|}{s} ds \leq \frac{1}{2} L M_\delta(u) u^2.$$

517 Thus, since parallel transport preserves inner products,

$$\begin{aligned}
 518 \quad & \frac{\|\dot{X}_\delta(t)\|}{t} \leq \left(\frac{1}{2} C L M_\delta(\delta) \delta^p + C \delta^p \|\text{gradf}(x_0)\| \right) \frac{e^{-\frac{v}{\delta}t}}{t} \int_0^t e^{\frac{v}{\delta}u} du \\
 519 \quad & \leq \left(\frac{1}{2} C L M_\delta(\delta) \delta^p + C \delta^p \|\text{gradf}(x_0)\| \right) \frac{\delta}{vt} (1 - e^{-\frac{v}{\delta}t}) \leq \frac{1}{2} C L M_\delta(\delta) \delta^p + C \delta^p \|\text{gradf}(x_0)\|.
 \end{aligned}$$

520 Taking the supremum over $0 < t \leq \delta$ and rearranging gives for $\delta < \delta_M = \left(\frac{2}{CL} \right)^{\frac{1}{p}}$ that

$$521 \quad M_\delta(\delta) \leq \frac{2C\delta^p \|\text{gradf}(x_0)\|}{2 - CL\delta^p}.$$

522 The case $p < 2$ is done exactly in the same way except that we do not need to bound u^{p-2} by δ^{p-2}
 523 in the integrals since the t^{p-2} term in the differential equation is already replaced by δ^{p-2} .

524 Note that when $\delta < \delta_M$ and $\delta < t < t_M = \left(\frac{2(v+p+1)}{CL} \right)^{\frac{1}{p}}$, the smoothed ODE can be rewritten as

$$525 \quad \frac{d}{dt} (t^v \dot{X}_\delta(t)) = -Ct^{v+p-2} \text{gradf}(X_\delta).$$

Therefore, we can use Lemma 4 in [4] once again to obtain

$$\begin{aligned} \Gamma_{X_\delta(t)}^{X_\delta(\delta)} t^v \dot{X}_\delta(t) - \delta^v \dot{X}_\delta(\delta) &= \int_0^t \left(\Gamma_{X_\delta(u)}^{X_\delta(\delta)} \text{gradf}(X_\delta(u)) - \Gamma_{X_\delta(u)}^{X_\delta(\delta)} \Gamma(X_\delta)_{x_0}^{X_\delta(u)} \text{gradf}(x_0) \right) C u^{v+p-2} du \\ &\quad - \int_0^t C u^{v+p-2} \Gamma_{X_\delta(u)}^{X_\delta(\delta)} \Gamma(X_\delta)_{x_0}^{X_\delta(u)} \text{gradf}(x_0) du. \end{aligned}$$

Using the fact that parallel transport preserves inner products, and dividing by t^{v+1} gives

$$\begin{aligned} \frac{\|\dot{X}_\delta(t)\|}{t} &\leq \frac{\delta^{v+1}}{t^{v+1}} \frac{\|\dot{X}_\delta(\delta)\|}{\delta} + \frac{CL}{2t^{v+1}} \int_\delta^t M_\delta(u) u^{v+p} du + \frac{C}{t^{v+1}} \|\text{gradf}(x_0)\| \int_\delta^t u^{v+p-2} du \\ &\leq \frac{\delta^{v+1}}{t^{v+1}} \frac{2C\delta^p \|\text{gradf}(x_0)\|}{2-CL\delta^p} + \frac{CL}{2(v+p+1)} M_\delta(t) t^p + \frac{C(t^{v+p-1} - \delta^{v+p-1})}{(v+p-1)t^{v+1}} \|\text{gradf}(x_0)\|, \end{aligned}$$

and since this upper bound is an increasing function of t , we have for any $t' \in (\delta, t)$ that

$$\frac{\|\dot{X}_\delta(t')\|}{t'} \leq \frac{2C\delta^p \|\text{gradf}(x_0)\|}{2-CL\delta^p} + \frac{CL}{2(v+p+1)} M_\delta(t) t^p + \frac{Ct^{p-2}}{v+p-1} \|\text{gradf}(x_0)\|.$$

Taking the supremum over all $t' \in (0, t)$ gives for $\delta < \delta_M$ and $\delta < t < t_M$,

$$M_\delta(t) \leq \frac{1}{1 - \frac{CL}{2(v+p+1)} t^p} \left(\frac{2C\delta^p}{2-CL\delta^p} + \frac{Ct^{p-2}}{v+p-1} \right) \|\text{gradf}(x_0)\|.$$

Now consider the family of functions

$$\mathcal{F} = \left\{ X_\delta : [0, T] \rightarrow \mathbb{R} \mid \delta = 2^{-n} \tilde{\delta}, n = 0, 1, \dots \right\},$$

where $T = \left(\frac{v+p+1}{CL}\right)^{\frac{1}{p}}$ and $\tilde{\delta} = \left(\frac{1}{CL}\right)^{\frac{1}{p}}$. By definition of M_δ , we have for $t \in [0, T]$ and $\delta \in (0, \tilde{\delta})$ that

$$\|\dot{X}_\delta\| \leq TM_\delta(T) \leq 2CT \left(\tilde{\delta} + \frac{CT^{p-2}}{v+p-1} \right) \quad \text{and} \quad d(X_\delta(t), X_\delta(0)) \leq \int_0^t \|\dot{X}_\delta(u)\| du \leq t \|\dot{X}_\delta\| \leq T \|\dot{X}_\delta\|.$$

Thus, \mathcal{F} is equicontinuous and uniformly bounded, and the Riemannian manifold \mathcal{Q} is complete by Assumption 1, so by the Arzela–Ascoli Theorem (Theorem 17 in [10]), \mathcal{F} contains a subsequence that converges uniformly on $[0, T]$ to some function X^* . The same argument as in part 5 of the proof of Lemma 3 of [4] shows that X^* is a solution to the original initial-valued ODE on $[0, T]$ which can then be extended to get a global solution on $[0, \infty)$. \square

B.2. Strongly Convex Case.

Theorem B.2. *Suppose that Assumption 1 is satisfied, and that $\eta > 0$ is a given constant. Then, the differential equation*

$$\nabla_{\dot{X}} \dot{X} + \eta \dot{X} + \text{gradf}(X) = 0,$$

has a global solution $X : [0, \infty) \rightarrow \mathcal{Q}$ under the initial conditions $X(0) = x_0 \in \mathcal{Q}$ and $\dot{X}(0) = 0$.

Proof. Exp and Log are defined globally on \mathcal{Q} by Assumption 1, so we can choose geodesically normal coordinates $\phi = \psi^{-1}$ around x_0 defined globally on \mathcal{Q} and put $c = \phi \circ X$. As in [4], using the smoothness of f and letting $u = \dot{c}$ gives a system of first order ODEs which defines a local representation for a vector field in $T\mathcal{Q}$, and results from section IV.3 of [11] guarantee that the initial-valued differential equation has a unique solution locally around 0. It remains to show that this solution actually exists on $[0, \infty)$. Towards contradiction, suppose $[0, T)$ is the maximal interval of existence of the solution X , for some finite $T > 0$. Then,

$$\frac{d}{dt} f(X(t)) = \langle \text{gradf}(X), \dot{X} \rangle = -\langle \nabla_{\dot{X}} \dot{X}, \dot{X} \rangle - C \langle \dot{X}, \dot{X} \rangle = -\frac{1}{2} \frac{d}{dt} \|\dot{X}\|^2 - C \|\dot{X}\|^2.$$

559 Rearranging, integrating both sides and using the Cauchy-Schwarz inequality gives

$$560 \int_0^T \|\dot{X}\| dt = \sqrt{T(f(x_0) - \inf_u f(u)) + \frac{T}{2} \left(\|\dot{X}(0)\|^2 - \inf_{t \in [0, T]} \|\dot{X}(t)\|^2 \right)} < \infty,$$

561 since f is bounded from below by Assumption 1. Thus, $\lim_{t \rightarrow T} X(t)$ exists, and since \mathcal{Q} is complete,
562 the limit is in \mathcal{Q} , contradicting the maximality of $[0, T)$, thereby concluding the proof. \square

563

564

APPENDIX C. PROOFS OF CONVERGENCE RATES

565 The proofs of the convergence rates of solutions to p -Bregman Euler–Lagrange equations are
566 inspired by those of Theorems 5 and 6 from [4], and make use of Lemmas 2 and 12 therein:

567 **Lemma C.1.** *Given a Riemannian manifold \mathcal{Q} with sectional curvature bounded above by K_{\max}
568 and below by K_{\min} , with ζ given by equation (2.1), and such that*

$$569 \text{diam}(\mathcal{Q}) < \begin{cases} \frac{\pi}{\sqrt{K_{\max}}} & \text{if } K_{\max} > 0 \\ \infty & \text{if } K_{\max} \leq 0 \end{cases},$$

570 we have that

$$571 \langle \nabla_{\dot{X}} \text{Log}_X(p), -\dot{X} \rangle \leq \zeta \|\dot{X}\|^2.$$

572 **Lemma C.2.** *Given a point q and a smooth curve $X(t)$ on a Riemannian manifold \mathcal{Q} ,*

$$573 \frac{d}{dt} \|\text{Log}_{X(t)}(q)\|^2 = 2 \langle \text{Log}_{X(t)}(q), \nabla_{\dot{X}} \text{Log}_{X(t)}(q) \rangle = 2 \langle \text{Log}_{X(t)}(q), -\dot{X}(t) \rangle.$$

574

C.1. Convex Case.

576 **Theorem C.1.** *Suppose $f : \mathcal{Q} \rightarrow \mathbb{R}$ is a geodesically convex function, and Assumption 1 is satisfied.
577 Then, any solution $X(t)$ of the p -Bregman Euler–Lagrange equation*

$$578 \nabla_{\dot{X}} \dot{X} + \frac{\zeta p + 1}{t} \dot{X} + Cp^2 t^{p-2} \text{grad}f(X) = 0,$$

579 converges to a minimizer x^* of f with rate

$$580 f(X(t)) - f(x^*) \leq \frac{\zeta \|\text{Log}_{x_0}(x^*)\|^2}{2Ct^p}.$$

581 *Proof.* Let

$$582 \mathcal{E}(t) = Ct^p (f(X) - f(x^*)) + \frac{1}{2} (\zeta - 1) \|\text{Log}_X(x^*)\|^2 + \frac{1}{2} \left\| \frac{t}{p} \dot{X} - \text{Log}_X(x^*) \right\|^2.$$

583 Then, using Lemma C.2,

$$\begin{aligned} 584 \dot{\mathcal{E}}(t) &= Cpt^{p-1} (f(X) - f(x^*)) + Ct^p \langle \text{grad}f(X), \dot{X} \rangle + (\zeta - 1) \langle \text{Log}_X(x^*), -\dot{X} \rangle \\ 585 &\quad + \left\langle \frac{t}{p} \dot{X} - \text{Log}_X(x^*), \frac{1}{p} \dot{X} + \frac{t}{p} \nabla_{\dot{X}} \dot{X} - \nabla_{\dot{X}} \text{Log}_X(x^*) \right\rangle \\ 586 &= Cpt^{p-1} (f(X) - f(x^*)) + Ct^p \langle \text{grad}f(X), \dot{X} \rangle + (\zeta - 1) \langle \text{Log}_X(x^*), -\dot{X} \rangle \\ 587 &\quad + \left\langle \frac{t}{p} \dot{X} - \text{Log}_X(x^*), \left(\frac{1}{p} \dot{X} + \frac{t}{p} \nabla_{\dot{X}} \dot{X} + \zeta \dot{X} \right) - \zeta \dot{X} - \nabla_{\dot{X}} \text{Log}_X(x^*) \right\rangle. \end{aligned}$$

588 Now, the p -Bregman Euler–Lagrange equation can be rewritten as

$$589 \frac{1}{p} \dot{X} + \frac{t}{p} \nabla_{\dot{X}} \dot{X} + \zeta \dot{X} = -Cpt^{p-1} \text{grad}f(X).$$

590 Thus,

$$591 \quad \dot{\mathcal{E}}(t) = Cpt^{p-1} (f(X) - f(x^*)) + Ct^p \langle \text{grad}f(X), \dot{X} \rangle + (\zeta - 1) \langle \text{Log}_X(x^*), -\dot{X} \rangle \\ 592 \quad \quad \quad + \left\langle \frac{t}{p} \dot{X} - \text{Log}_X(x^*), -Cpt^{p-1} \text{grad}f(X) - \zeta \dot{X} - \nabla_{\dot{X}} \text{Log}_X(x^*) \right\rangle.$$

593 Canceling the $\langle \text{grad}f(X), \dot{X} \rangle$ and $\langle \text{Log}_X(x^*), -\dot{X} \rangle$ terms out using Lemma C.2, we get

$$594 \quad \dot{\mathcal{E}}(t) = Cpt^{p-1} [f(X) - f(x^*) + \langle \text{Log}_X(x^*), \text{grad}f(X) \rangle] - \frac{t}{p} (\langle \dot{X}, \nabla_{\dot{X}} \text{Log}_X(x^*) \rangle + \zeta \langle \dot{X}, \dot{X} \rangle).$$

595 Now, since f is geodesically convex, we have that $[f(X) - f(x^*) + \langle \text{Log}_X(x^*), \text{grad}f(X) \rangle] \leq 0$.

596 Furthermore, Lemma C.1 ensures that $(\langle \dot{X}, \nabla_{\dot{X}} \text{Log}_X(x^*) \rangle + \zeta \langle \dot{X}, \dot{X} \rangle) \geq 0$. Thus, $\dot{\mathcal{E}}(t) \leq 0$, so

$$597 \quad Ct^p (f(X) - f(x^*)) \leq Ct^p (f(X) - f(x^*)) + \frac{1}{2} (\zeta - 1) \|\text{Log}_X(x^*)\|^2 + \frac{1}{2} \left\| \frac{t}{p} \dot{X} - \text{Log}_X(x^*) \right\|^2 \\ 598 \quad \quad \quad = \mathcal{E}(t) \leq \mathcal{E}(0) = \frac{1}{2} (\zeta - 1) \|\text{Log}_{x_0}(x^*)\|^2 + \frac{1}{2} \|\text{Log}_{x_0}(x^*)\|^2 = \frac{1}{2} \zeta \|\text{Log}_{x_0}(x^*)\|^2.$$

599 which gives the desired rate of convergence. \square

600

601 C.2. Weakly-Quasi-Convex Case.

602 **Theorem C.2.** *Suppose $f : \mathcal{Q} \rightarrow \mathbb{R}$ is a geodesically α -weakly-quasi-convex function, and suppose*
603 *that Assumption 1 is satisfied. Then, any solution $X(t)$ of the p -Bregman Euler–Lagrange equation*

$$604 \quad \nabla_{\dot{X}} \dot{X} + \frac{\zeta p + \alpha}{\alpha t} \dot{X} + Cp^2 t^{p-2} \text{grad}f(X) = 0,$$

605 *converges to a minimizer x^* of f with rate*

$$606 \quad f(X(t)) - f(x^*) \leq \frac{\zeta \|\text{Log}_{x_0}(x^*)\|^2}{2C\alpha^2 t^p}.$$

607 *Proof.* Let

$$608 \quad \mathcal{E}(t) = C\alpha^2 t^p (f(X) - f(x^*)) + \frac{1}{2} (\zeta - 1) \|\text{Log}_X(x^*)\|^2 + \frac{1}{2} \left\| \frac{\alpha t}{p} \dot{X} - \text{Log}_X(x^*) \right\|^2.$$

609 Then, using Lemma C.2,

$$610 \quad \dot{\mathcal{E}}(t) = Cp\alpha^2 t^{p-1} (f(X) - f(x^*)) + C\alpha^2 t^p \langle \text{grad}f(X), \dot{X} \rangle + (\zeta - 1) \langle \text{Log}_X(x^*), -\dot{X} \rangle \\ 611 \quad \quad \quad + \left\langle \frac{\alpha t}{p} \dot{X} - \text{Log}_X(x^*), \frac{\alpha}{p} \dot{X} + \frac{\alpha t}{p} \nabla_{\dot{X}} \dot{X} - \nabla_{\dot{X}} \text{Log}_X(x^*) \right\rangle \\ 612 \quad \quad \quad = Cp\alpha^2 t^{p-1} (f(X) - f(x^*)) + C\alpha^2 t^p \langle \text{grad}f(X), \dot{X} \rangle + (\zeta - 1) \langle \text{Log}_X(x^*), -\dot{X} \rangle \\ 613 \quad \quad \quad + \left\langle \frac{\alpha t}{p} \dot{X} - \text{Log}_X(x^*), \left(\frac{\alpha}{p} \dot{X} + \frac{\alpha t}{p} \nabla_{\dot{X}} \dot{X} + \zeta \dot{X} \right) - \zeta \dot{X} - \nabla_{\dot{X}} \text{Log}_X(x^*) \right\rangle.$$

614 Now, the p -Bregman Euler–Lagrange equation can be rewritten as

$$615 \quad \frac{\alpha}{p} \dot{X} + \frac{\alpha t}{p} \nabla_{\dot{X}} \dot{X} + \zeta \dot{X} = -Cp\alpha t^{p-1} \text{grad}f(X).$$

616 Thus,

$$617 \quad \dot{\mathcal{E}}(t) = Cp\alpha^2 t^{p-1} (f(X) - f(x^*)) + C\alpha^2 t^p \langle \text{grad}f(X), \dot{X} \rangle + (\zeta - 1) \langle \text{Log}_X(x^*), -\dot{X} \rangle \\ 618 \quad \quad \quad + \left\langle \frac{\alpha t}{p} \dot{X} - \text{Log}_X(x^*), -Cp\alpha t^{p-1} \text{grad}f(X) - \zeta \dot{X} - \nabla_{\dot{X}} \text{Log}_X(x^*) \right\rangle.$$

619 Canceling the $\langle \text{grad}f(X), \dot{X} \rangle$ and $\langle \text{Log}_X(x^*), -\dot{X} \rangle$ terms out using Lemma C.2, we get

$$620 \quad \dot{\mathcal{E}}(t) = Cp\alpha t^{p-1} [\alpha (f(X) - f(x^*)) + \langle \text{Log}_X(x^*), \text{grad}f(X) \rangle] - \frac{\alpha t}{p} (\langle \dot{X}, \nabla_{\dot{X}} \text{Log}_X(x^*) \rangle + \zeta \langle \dot{X}, \dot{X} \rangle).$$

621 Now, f is geodesically α -weakly-quasi-convex, so $[\alpha (f(X) - f(x^*)) + \langle \text{Log}_X(x^*), \text{grad}f(X) \rangle] \leq 0$.
622 Furthermore, Lemma C.1 ensures that $(\langle \dot{X}, \nabla_{\dot{X}} \text{Log}_X(x^*) \rangle + \zeta \langle \dot{X}, \dot{X} \rangle) \geq 0$. Thus, $\dot{\mathcal{E}}(t) \leq 0$, so

$$623 \quad C\alpha^2 t^p (f(X) - f(x^*)) \leq C\alpha^2 t^p (f(X) - f(x^*)) + \frac{1}{2}(\zeta - 1) \|\text{Log}_X(x^*)\|^2 + \frac{1}{2} \left\| \frac{\alpha t}{p} \dot{X} - \text{Log}_X(x^*) \right\|^2$$

$$624 \quad = \mathcal{E}(t) \leq \mathcal{E}(0) = \frac{1}{2}(\zeta - 1) \|\text{Log}_{x_0}(x^*)\|^2 + \frac{1}{2} \|\text{Log}_{x_0}(x^*)\|^2 = \frac{1}{2} \zeta \|\text{Log}_{x_0}(x^*)\|^2,$$

625 which gives the desired rate of convergence. \square

626

627

APPENDIX D. PROOF OF INVARIANCE THEOREM

628 **Theorem D.1.** *Suppose Assumption 1 is satisfied and that the curve $X(t)$ satisfies a p -Bregman*
629 *Euler–Lagrange equation of the form*

$$630 \quad \nabla_{\dot{X}} \dot{X} + \frac{\lambda p + 1}{t} \dot{X} + Cp^2 t^{p-2} \text{grad}f(X) = 0,$$

631 *for some $\lambda \in \mathbb{R}$. Then the reparametrized curve $X(t^{\hat{p}/p})$ satisfies the corresponding \hat{p} -Bregman*
632 *Euler–Lagrange equation.*

633 *Proof.* Let $\tau(t) = t^{\hat{p}/p}$ and $Y(t) = X(\tau(t))$. Then

$$634 \quad \dot{Y}(t) = \dot{\tau}(t) \dot{X}(\tau(t)), \quad \text{and} \quad \nabla_{\dot{Y}(t)} \dot{Y}(t) = \ddot{\tau}(t) \dot{X}(\tau(t)) + \dot{\tau}^2(t) \nabla_{\dot{X}(\tau(t))} \dot{X}(\tau(t)).$$

635 Inverting these relations gives

$$636 \quad \dot{X}(\tau(t)) = \frac{1}{\dot{\tau}(t)} \dot{Y}(t), \quad \text{and} \quad \nabla_{\dot{X}(\tau(t))} \dot{X}(\tau(t)) = \frac{1}{\dot{\tau}^2(t)} \nabla_{\dot{Y}(t)} \dot{Y}(t) - \frac{\ddot{\tau}(t)}{\dot{\tau}^3(t)} \dot{Y}(t).$$

637 The p -Bregman Euler–Lagrange equation at time $\tau(t)$ is given by

$$638 \quad \nabla_{\dot{X}(\tau(t))} \dot{X}(\tau(t)) + \frac{\lambda p + 1}{\tau(t)} \dot{X}(\tau(t)) + Cp^2 \tau^{p-2}(t) \text{grad}f(X(\tau(t))) = 0.$$

639 Substituting the expressions for $X(\tau(t))$, $\dot{X}(\tau(t))$ and $\nabla_{\dot{X}(\tau(t))} \dot{X}(\tau(t))$ in terms of $Y(t)$ and its
640 derivatives, and multiplying by $\dot{\tau}^2(t)$ gives

$$641 \quad \nabla_{\dot{Y}(t)} \dot{Y}(t) + \left((\lambda p + 1) \frac{\dot{\tau}(t)}{\tau(t)} - \frac{\ddot{\tau}(t)}{\dot{\tau}(t)} \right) \dot{Y}(t) + Cp^2 \dot{\tau}^2(t) \tau^{p-2}(t) \text{grad}f(Y(t)) = 0.$$

642 Substituting $\tau(t) = t^{\frac{\hat{p}}{p}}$ yields the \hat{p} -Bregman Euler–Lagrange equation for Y at time t . \square

643

644

REFERENCES

- 645 [1] P.-A. Absil, R. Mahony, and R. Sepulchre. *Optimization Algorithms on Matrix Manifolds*,
646 volume 78. 12 2008. ISBN 978-0-691-13298-3. doi: 10.1515/9781400830244.
- 647 [2] K. Ahn and S. Sra. From nesterov’s estimate sequence to riemannian acceleration, 2020.
- 648 [3] F. Alimisis, A. Orvieto, G. Bécigneul, and A. Lucchi. Practical accelerated optimization on
649 Riemannian manifolds, 2020.
- 650 [4] F. Alimisis, A. Orvieto, G. Bécigneul, and A. Lucchi. A continuous-time perspective for
651 modeling acceleration in Riemannian optimization. In *Proceedings of the 23rd International*
652 *AISTATS Conference*, volume 108 of *PMLR*, pages 1297–1307, 2020.

- 653 [5] M. Betancourt, M. I. Jordan, and A. Wilson. On symplectic optimization. 2018.
- 654 [6] A.-L. Cauchy. Méthode générale pour la résolution des systèmes d'équations simultanées. *Acad.*
655 *Sci. Paris*, 25:536–538, 1847.
- 656 [7] V. Duruisseaux, J. Schmitt, and M. Leok. Adaptive Hamiltonian variational integrators and
657 applications to symplectic accelerated optimization, 2020.
- 658 [8] E. Hairer, C. Lubich, and G. Wanner. *Geometric numerical integration*, volume 31 of *Springer*
659 *Series in Computational Mathematics*. Springer-Verlag, Berlin, second edition, 2006.
- 660 [9] J. Jost. *Riemannian geometry and geometric analysis*. Universitext. Springer, Cham, seventh
661 edition, 2017.
- 662 [10] J. L. Kelley. *General Topology*. Graduate Texts in Mathematics. Springer New York, 1975.
663 ISBN 9780387901251.
- 664 [11] S. Lang. *Fundamentals of Differential Geometry*, volume 191 of *Graduate Texts in Mathemat-*
665 *ics*. Springer -Verlag, New York, 1999. ISBN 9780387985930.
- 666 [12] J. M. Lee. *Introduction to Riemannian Manifolds*, volume 176 of *Graduate Texts in Mathe-*
667 *matics*. Springer, Cham, second edition, 2018.
- 668 [13] M. Leok and T. Ohsawa. Variational and geometric structures of discrete Dirac mechanics.
669 *Found. Comput. Math.*, 11(5):529–562, 2011.
- 670 [14] Y. Liu, F. Shang, J. Cheng, H. Cheng, and L. Jiao. Accelerated first-order methods for
671 geodesically convex optimization on Riemannian manifolds. In *Advances in Neural Information*
672 *Processing Systems*, volume 30, pages 4868–4877. Curran Associates, Inc., 2017.
- 673 [15] J. E. Marsden and T. S. Ratiu. *Introduction to mechanics and symmetry*, volume 17 of *Texts*
674 *in Applied Mathematics*. Springer-Verlag, New York, second edition, 1999.
- 675 [16] J. E. Marsden and M. West. Discrete mechanics and variational integrators. *Acta Numer.*, 10:
676 357–514, 2001.
- 677 [17] A. S. Nemirovsky and D. B. Yudin. *Problem Complexity and Method Efficiency in Optimiza-*
678 *tion*. Wiley - Interscience series in discrete mathematics. Wiley, 1983.
- 679 [18] Y. Nesterov. A method of solving a convex programming problem with convergence rate
680 $\mathcal{O}(1/k^2)$. *Soviet Mathematics Doklady*, 27(2):372–376, 1983.
- 681 [19] Y. Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*, volume 87 of
682 *Applied Optimization*. Kluwer Academic Publishers, Boston, MA, 2004.
- 683 [20] Y. Nesterov. Accelerating the cubic regularization of Newton's method on convex problems.
684 *Math. Program.*, 112:159–181, 2008.
- 685 [21] A. Orvieto and A. Lucchi. Shadowing properties of optimization algorithms. In *Advances in*
686 *Neural Information Processing Systems*, volume 32, pages 12692–12703, 2019.
- 687 [22] W. Su, S. Boyd, and E. Candes. A differential equation for modeling Nesterov's Accelerated
688 Gradient method: theory and insights. *Journal of Machine Learning Research*, 17(153):1–43,
689 2016.
- 690 [23] I. Sutskever, J. Martens, G. Dahl, and G. Hinton. On the importance of initialization and mo-
691 mentum in deep learning. In *Proceedings of the 30th International Conference on International*
692 *Conference on Machine Learning - Volume 28*, ICML'13, page III–1139–III–1147, 2013.
- 693 [24] A. Wibisono, A. Wilson, and M. Jordan. A variational perspective on accelerated methods in
694 optimization. *Proceedings of the National Academy of Sciences*, 113(47):E7351–E7358, 2016.
- 695 [25] H. Zhang and S. Sra. First-order methods for geodesically convex optimization. In *29th Annual*
696 *Conference on Learning Theory*, pages 1617–1638, 2016.
- 697 [26] H. Zhang and S. Sra. Towards riemannian accelerated gradient methods. 2018.
- 698 [27] J. Zhang, A. Mokhtari, S. Sra, and A. Jadbabaie. Direct runge-kutta discretization achieves
699 acceleration. 2018.