

# On Symplectic Optimization

Michael Betancourt  
Michael I. Jordan  
Ashia C. Wilson

*Abstract.* Accelerated gradient methods have had significant impact in machine learning—in particular the theoretical side of machine learning—due to their ability to achieve oracle lower bounds. But their heuristic construction has hindered their full integration into the practical machine-learning algorithmic toolbox, and has limited their scope. In this paper we build on recent work which casts acceleration as a phenomenon best explained in continuous time, and we augment that picture by providing a systematic methodology for converting continuous-time dynamics into discrete-time algorithms while retaining oracle rates. Our framework is based on ideas from Hamiltonian dynamical systems and symplectic integration. These ideas have had major impact in many areas in applied mathematics, but have not yet been seen to have a relationship with optimization.

## 1. INTRODUCTION

Optimization theory has played an increasingly central role in the development of machine learning in recent years. This has happened not only because optimization theory supplies algorithms and convergence rates for learning algorithms, but also because it supplies lower bounds, and hence fundamental understanding. A milestone in this regard was the discovery by [Nemirovskii & Yudin \(1983\)](#) of oracle lower bounds for gradient-based optimization, and the ensuing derivation by [Nesterov \(1983\)](#) of an “accelerated gradient descent” (AGD) algorithm whose rate is provably better than that of gradient descent, and which matches the oracle lower bound.

A flurry of mathematical and algorithmic results have followed in the wake of these seminal discoveries from the 1980’s, but even after three decades there remains a lack of

---

*Michael Betancourt is a research scientist at Symplectomorphic, LLC (e-mail: [betanalpha@gmail.com](mailto:betanalpha@gmail.com)). Michael I. Jordan is the Pehong Chen Distinguished Professor in the Departments of EECS and Statistics at the University of California, Berkeley (e-mail: [jordan@cs.berkeley.edu](mailto:jordan@cs.berkeley.edu)). Ashia C. Wilson is a graduate student in the Department of Statistics at the University of California, Berkeley.*

understanding of the general acceleration phenomenon. In particular, a theoretical framework that can *generate* accelerated methods has not yet emerged. Recent progress in this regard has been achieved by considering continuous-time analogs of acceleration methods (Su et al., 2016; Krichene et al., 2015). Notably, Wibisono et al. (2016) presented a variational framework, involving a “Bregman Lagrangian,” that generates differential equations with rates that are continuous-time analogs of the discrete-time oracle rates.

The work of Wibisono et al. (2016), however, only partially addresses the problem of providing a generative framework for acceleration. They show that any desired rate can be achieved in continuous time—different algorithms follow the same *path* in phase space while doing so at different *speeds*, and that the speed can be arbitrarily fast. Differences in speed thus correspond to a mere change of the clock by which time is measured. The fact that lower bounds for rates emerge in discrete time must therefore have something to do with the discretization of the class of differential equations arising from the Bregman Lagrangian. Wibisono et al. (2016) were able to provide an adhoc discretization that yielded an algorithm whose rate matches the rate of Nesterov acceleration in a particular setting, but their framework is silent on a general methodology for providing such discretizations.

The class of differential equations arising from the Bregman Lagrangian must be special in some sense, given that they deliver continuous-time analogs of oracle rates. The notion that certain differential equations are special has a long history in physics, where underlying Lagrangians and Hamiltonians possess certain mathematical symmetries that yield conservation laws for the resulting differential equations. Moreover, the venerable field of symplectic integration shows that it is possible to preserve these conservation laws when discretizing the differential equations (Hairer et al., 2006). The resulting integrators improve upon classical integrators, e.g., Euler, Runge-Kutta, precisely because they respect the underlying mathematical symmetries. This results in certain error terms canceling, with a variety of favorable consequences, including long-term stability. Of particular relevance to the current setting, the stability of such integrators means that it is possible to take much larger step sizes than with classical integrators. Given that we are interested in “accelerated” methods that arrive at an optimum as quickly as possibly, this feature of symplectic integration seems directly relevant.

In the current paper we show how to apply symplectic integration to gradient-based optimization. Our approach is cast in a Hamiltonian framework, obtained from the Bregman-Lagrangian framework via a Legendre transformation. This Hamiltonian is time-varying, a fact that we address via a lifting procedure. We then show how to derive a symplectic integrator from the lifted Hamiltonian. The end result is a fully generative mathematical pipeline, from problem specification to discrete-time accelerated algorithm. Our algorithms are related to Nesterov’s algorithms, but they are *not* exactly the same, and the differences are interesting.

## 2. BREGMAN DYNAMICS FOR OPTIMIZATION

We begin with a brief review of the dynamical framework introduced in [Wibisono et al. \(2016\)](#), including the heuristic discretization that those authors employed to obtain accelerated discrete-time optimization algorithms.

Consider a Euclidean vector space,  $\mathcal{X}$ , which we will denote as the *configuration manifold*. This configuration manifold is equipped with the Euclidean gradient operator,  $\nabla$ , and inner product,  $\langle \cdot, \cdot \rangle$ . Formally we should be careful to distinguish between points, vectors, and covectors on the configuration manifold but we will reserve that level of rigor for a more formal geometric treatment presented in [Appendix A](#).

Given a smooth *objective function* on the configuration manifold,  $f : \mathcal{X} \rightarrow \mathbb{R}$ , the optimization problem is then to compute minima of  $f$ :

$$x^* \in \underset{x \in \mathcal{X}}{\operatorname{argmin}} f(x).$$

In accordance with most of the theoretical literature on acceleration, we will focus on the setting in which  $f$  is convex and exhibits a single minimum. But it is worth emphasizing that convexity is not needed in our construction of the Hamiltonian nor for the symplectic integrators. Note, moreover, that there is a growing literature on acceleration in the non-convex setting, where the dynamical systems perspective presented here is also proving to be useful; see, e.g., [Jin et al. \(2017\)](#).

From a dynamical perspective, the objective function naturally plays the role of a potential energy, with the minimum at the basin of that potential. If we want to generate dynamics that might settle into this basin, however, then we need to consider not the configuration manifold but rather its *tangent bundle*,  $T\mathcal{X}$ , consisting of points in  $\mathcal{X}$  paired with tangent vectors or *velocities*. In particular, we need to complement the potential energy with a kinetic energy function on the tangent bundle.

Following [Wibisono et al. \(2016\)](#) we construct a kinetic energy from an auxiliary smooth function,  $h : \mathcal{X} \rightarrow \mathbb{R}$ , and its associated *Bregman divergence*,  $D_h(y, x) = h(y) - h(x) - \langle \nabla h(x), y - x \rangle$ . For any given point in the tangent bundle,  $(x, v) \in T\mathcal{X}$ , we can translate the base point,  $x$ , in the direction of the velocity,  $v$ , to give the new point  $x' = x + e^{-\alpha(t)}v$ , for a scaling function  $\alpha(t)$ . The divergence between these two points defines the *Bregman kinetic energy*:

$$\begin{aligned} K(x, v) &\equiv D_h(x + e^{-\alpha(t)}v, x) \\ &= h(x + e^{-\alpha(t)}v) - h(x) - e^{-\alpha(t)} \langle \nabla h(x), v \rangle. \end{aligned}$$

This kinetic energy admits the evocative interpretation as a comparison of how  $h$  changes under a finite translation:

$$\Delta h(x, v, t) = h(x + e^{-\alpha(t)}v) - h(x),$$

versus a scaled infinitesimal translation:

$$\delta h(x, v, t) = e^{-\alpha(t)} \langle \nabla h(x), v \rangle.$$

Defining a time-dependent potential energy,

$$U(x, t) = e^{\beta(t)} f(x),$$

we can then construct the *Bregman Lagrangian* as:

$$\begin{aligned} L(x, v, t) &= e^{\alpha(t)+\gamma(t)} (K(x, v) - U(x)) \\ &= e^{\alpha(t)+\gamma(t)} \left( D_h(x + e^{-\alpha(t)} v, x) - e^{\beta(t)} f(x) \right). \end{aligned}$$

From the Bregman Lagrangian we obtain a variational problem on the tangent bundle whose solutions yield smooth trajectories via the ordinary differential equations

$$\frac{d}{dt} \frac{\partial L}{\partial v}(x, v, t) = \frac{\partial L}{\partial x}(x, v, t).$$

The time dependence of the Lagrangian allows the dynamics to rapidly converge to a minimum, as opposed to dynamics obtained from a time-independent Lagrangian which would oscillate around the desired minimum.

[Wibisono et al. \(2016\)](#) also defined the following *ideal scaling conditions*:

$$\begin{aligned} \alpha(t) &= \log p - \log t \\ \beta(t) &= p \log t + \log C \\ \gamma(t) &= p \log t, \end{aligned}$$

for  $p, C \in \mathbb{R}^+$ , and demonstrated that if  $f$  and  $h$  are sufficiently well behaved then the Bregman dynamics will provably converge to the minimum of  $f$  at the polynomial rate,  $\mathcal{O}(1/t^p)$ . This captures not only classical Nesterov acceleration, with its rate of  $\mathcal{O}(1/t^2)$ , but also higher-order accelerated algorithms for which  $p > 2$ .

Unfortunately it is not obvious how to discretize these continuous dynamics to obtain a discrete-time algorithm. [Wibisono et al. \(2016\)](#) found that simple discretizations yield algorithms that do not recover accelerated Nesterov methods and can even be unstable. Ultimately they were able to find a stable discretization that yielded the oracle rates. Their discretization was a sophisticated but heuristic discretization, coupling a Crank-Nicolson discretization of the position updates and a backwards Euler discretization of the velocity updates with a third implicit sequence of intermediate positions,  $(y_n)$ :

$$\begin{aligned} x_{n+1} &= \frac{p}{n+p} z_n + \frac{n}{n+p} y_n \\ y_n &= \operatorname{argmin}_{y \in \mathcal{X}} \left[ f_{p-1}(y, ; x_n) + \frac{N}{\epsilon^p p} \|y - x_n\|^p \right] \\ z_n &= \operatorname{argmin}_{z \in \mathcal{X}} \left[ C p n^{p-1} \langle \nabla f(y_n), z \rangle + \frac{1}{\epsilon^p} D_h(z, z_{n-1}) \right]. \end{aligned}$$

Here  $f_{p-1}(y; x_n)$  is the order- $p$  Taylor expansion of the objective function around  $x_n$ . This third sequence proved to be the key, stabilizing the discretization and preserving the convergence rates of the continuous-time dynamics. This discretization was obtained by analogy with Nesterov's classical updates, and in this paper we will refer to these discretizations as *generalized Nesterov discretizations*.

### 3. SIMULATING BREGMAN DYNAMICS WITH SYMPLECTIC INTEGRATORS

The difficulties associated with discretizing Lagrangian dynamics are well known in the mathematical literature (Leimkuhler & Reich, 2004; Hairer et al., 2006). Discretizations of Lagrangian dynamics are often fragile, especially when the dimension of the configuration space is large. Even high-order discretizations can diverge after short integration times. Ultimately this is because any discretization of dynamics on the tangent bundle does not preserve the continuous symmetries of the dynamical system that stabilize the exact dynamics.

Fortunately we can readily construct a discretization that *does* preserve the necessary symmetries, by exploiting the dual *Hamiltonian* representation of the Bregman dynamics. The Hamiltonian system is the Legendre transform of the Lagrangian system, trading velocities,  $v$ , and the tangent bundle,  $T\mathcal{X}$ , for dual *momenta*,  $r$ , and the *cotangent bundle*.

In this section we will first construct the Hamiltonian representation of the Bregman dynamics and then modify that system to circumvent the explicit time dependence and admit the application of symplectic integrators.

#### 3.1 The Bregman Hamiltonian

To build up the Bregman Hamiltonian from the Bregman Lagrangian we first have to relate define momenta as the derivative of the Lagrangian with respect to the velocities,

$$\begin{aligned} r(x, v, t) &= \frac{\partial L}{\partial v}(x, v, t) \\ &= e^{\gamma(t)} \left( \frac{\partial h}{\partial x}(x + e^{-\alpha(t)}v) - \frac{\partial h}{\partial x}(x) \right). \end{aligned}$$

Given the Legendre conjugate of  $h$ ,

$$h^* = \sup_{v \in T\mathcal{X}} [r \cdot v - h(v)],$$

we can invert this relationship to give

$$v(x, r, t) = e^{\alpha(t)} \left( \frac{\partial h^*}{\partial r}(e^{-\gamma(t)}r + \frac{\partial h}{\partial x}(x)) - r \right).$$

We are now in position to construct the *Bregman Hamiltonian*:

$$\begin{aligned} H(x, r, t) &= r \cdot v(x, r, t) - L(x, v(r), t) \\ &= e^{\alpha(t)+\gamma(t)} \left( D_{h^*}(e^{-\gamma(t)} r + \frac{\partial h}{\partial x}(x), \frac{\partial h}{\partial x}(x)) + e^{\beta(t)} f(x) \right) \\ &= e^{\alpha(t)+\gamma(t)} \left( D_{h^*}(e^{-\gamma(t)} r + \nabla h(x), \nabla h(x)) + e^{\beta(t)} f(x) \right), \end{aligned}$$

where

$$D_{h^*}(r, s) = h^*(r) - h^*(s) - \frac{\partial h^*}{\partial r}(s) \cdot (p - s).$$

The Bregman dynamics can then be generated from the Bregman Hamiltonian by integrating *Hamilton's equations*:

$$\frac{dx}{dt} = + \frac{\partial H}{\partial r}(x, r, t), \quad \frac{dr}{dt} = - \frac{\partial H}{\partial x}(x, r, t).$$

When the Hamiltonian does not explicitly depend on time the dynamics are said to be *autonomous* (José & Saletan, 1998) and there are standard methods for constructing symplectic integrators that preserve the critical symmetries that stabilize the dynamics. These integrators are extremely accurate, defining discretized dynamics that mirror the exact dynamics even for long integration times and high-dimensional configuration spaces (Leimkuhler & Reich, 2004; Hairer et al., 2006).

Unfortunately the explicit time dependence that allows the dynamics to converge to the minimum of the objective also renders the Bregman Hamiltonian *non-autonomous*. Fortunately this problem can be circumvented. As we show in the next section, by introducing a few more auxiliary variables we can lift the non-autonomous Bregman Hamiltonian system into an autonomous, *extended* Hamiltonian system where symplectic integrators are immediately applicable.

### 3.2 Making the non-autonomous autonomous

For an autonomous dynamical system time is simply a parameterization of motion along a dynamical trajectory. In particular, we can utilize uniform time increments to facilitate stable discretization of those dynamics. When the trajectories themselves explicitly depend on time, however, stable discretizations become all the more challenging. In order to overcome this difficulty we need to decouple the two responsibilities of time in a dynamical system by incorporating the explicit time into the configuration space and introducing a new effective time to parameterize motion along the dynamical trajectories (de León & Rodrigues, 1989).

The explicit time,  $t$ , serves as a new position in the extended configuration space,  $(x, t) \in \Xi$ . The extended cotangent bundle then includes a conjugate energy,  $(x, t, r, \mathcal{E}) \in T^*\Xi$ . We then define the extended Hamiltonian as

$$H_{\Xi} = \mathcal{E} - H(x, t, r).$$

Here the conjugate energy  $\mathcal{E}$  must compensate for the time dependence of the original Hamiltonian to ensure that the extended Hamiltonian,  $H_{\Xi}$ , is constant along dynamical trajectories.

The corresponding equations of motion for the extended Hamiltonian system become

$$\begin{aligned}\frac{dx}{d\tau} &= +\frac{\partial H_{\Xi}}{\partial r}(x, t, r, \mathcal{E}) \\ \frac{dt}{d\tau} &= +\frac{\partial H_{\Xi}}{\partial \mathcal{E}}(x, t, r, \mathcal{E}) \\ \frac{dr}{d\tau} &= -\frac{\partial H_{\Xi}}{\partial x}(x, t, r, \mathcal{E}) \\ \frac{d\mathcal{E}}{d\tau} &= -\frac{\partial H_{\Xi}}{\partial t}(x, t, r, \mathcal{E}),\end{aligned}$$

which introduces the effective time,  $\tau$ , to parameterize the motion along the extended dynamical trajectories. These dynamics projected back down to the original cotangent bundle yield the original Bregman dynamics, but by separating the time dependence of the dynamics from the parameterization of the trajectories these extended dynamics are manifestly autonomous. In particular, we can immediately apply symplectic integrators to the extended Hamiltonian system.

### 3.3 Building an extended leapfrog integrator

We are now in a position to build a symplectic integrator to simulate the Bregman dynamics. There is a rich literature on symplectic integrators and their exceptional performance (Leimkuhler & Reich, 2004; Hairer et al., 2006), so here we will limit our discussion to the construction and behavior of a simple *leapfrog integrator* for our extended Hamiltonian system. Despite the simplicity of this integrator, we will see in Section 4 that it rivals the performance of the generalized Nesterov discretization.

Symplectic integrators are naturally constructed by splitting the Hamiltonian into component Hamiltonians whose dynamics can be solved exactly, or at least sufficiently close to exactly numerically, and then composing those dynamics together symmetrically. For example, consider the splitting

$$H_{\Xi} = H_A + H_B + H_C,$$

where

$$\begin{aligned}H_A(\mathcal{E}) &= \mathcal{E} \\ H_B(x, r, t) &= e^{\alpha(t)+\gamma(t)} D_{h^*}(e^{-\gamma(t)}r + \frac{\partial h}{\partial x}(r), \frac{\partial h}{\partial x}(x)) \\ H_C(x, t) &= e^{\alpha(t)+\gamma(t)+\beta(t)} f(x).\end{aligned}$$

These three component Hamiltonians generate dynamics in the extended cotangent bundle with the six vector fields,

$$\begin{aligned}\vec{H}_A &= +\frac{\partial H_1}{\partial \mathcal{E}}(\mathcal{E}) \frac{d}{dt} \\ \vec{H}_{B1} &= -\frac{\partial H_2}{\partial x}(x, r, t) \frac{d}{dr} \\ \vec{H}_{B2} &= -\frac{\partial H_2}{\partial t}(x, r, t) \frac{d}{d\mathcal{E}} \\ \vec{H}_{B3} &= +\frac{\partial H_2}{\partial r}(x, r, t) \frac{d}{dx} \\ \vec{H}_{C1} &= -\frac{\partial H_3}{\partial x}(x, t) \frac{d}{dr} \\ \vec{H}_{C2} &= -\frac{\partial H_3}{\partial t}(x, t) \frac{d}{d\mathcal{E}}.\end{aligned}$$

Regardless of the nature of  $h$ , the vector fields  $\vec{H}_A$ ,  $\vec{H}_{B2}$ ,  $\vec{H}_{C1}$ , and  $\vec{H}_{C2}$  will always be trivial and hence their evolution can be solved exactly. For example,

$$\exp\left(\epsilon \vec{H}_{C1}\right)(x, t, r, \mathcal{E}) = \left(x, t, r - \epsilon \frac{\partial H_3}{\partial x}(x, t), \mathcal{E}\right).$$

On the other hand the component dynamics of  $\vec{H}_{B1}$  and  $\vec{H}_{B3}$  may or may not be trivial, depending on the choice of  $h$ . Even if they are nonlinear, however, they can be solved implicitly using, for example, fixed-point iterations.

We can now build a symplectic integrator by composing these component dynamics together to approximate the full dynamics. Here we will consider a symmetric *leapfrog* composition,

$$\begin{aligned}\Phi_\epsilon &= \exp\left(\frac{\epsilon}{2}\vec{H}_A\right) \circ \exp\left(\frac{\epsilon}{2}\left(\vec{H}_{B2} + \vec{H}_{C2}\right)\right) \circ \exp\left(\frac{\epsilon}{2}\vec{H}_{C1}\right) \\ &\quad \circ \exp\left(\frac{\epsilon}{2}\vec{H}_{B1}\right) \circ \exp\left(\epsilon\vec{H}_{B3}\right) \circ \exp\left(\frac{\epsilon}{2}\vec{H}_{B1}\right) \\ &\quad \circ \exp\left(\frac{\epsilon}{2}\vec{H}_{C1}\right) \circ \exp\left(\frac{\epsilon}{2}\left(\vec{H}_{B2} + \vec{H}_{C2}\right)\right) \circ \exp\left(\frac{\epsilon}{2}\vec{H}_A\right).\end{aligned}$$

Applying the Baker-Campbell-Hausdorff equation to this composite operator demonstrates that the symmetry of the composition ensures the cancellation of all terms linear and quadratic in the step size, leaving

$$\Phi_\epsilon = \exp\left(\epsilon \vec{H}\right) + \mathcal{O}(\epsilon^3).$$

If we apply the composite operator for  $N = T/\epsilon$  steps we then we can approximate the evolution of the exact dynamics for time  $T$  with error only quadratic in the step size:

$$(\Phi_\epsilon)^N = \left(\exp\left(\epsilon \vec{H}\right)\right)^N + \frac{T}{\epsilon} \mathcal{O}(\epsilon^3) = \exp\left(T \vec{H}\right) + \mathcal{O}(\epsilon^2).$$



Because each component operator exactly solves a dynamical system, their solutions preserve the dynamical symmetries that maintain stable evolution. Moreover, the symmetric composition of these component dynamics yields an approximate dynamics that very accurately tracks the dynamics of the extended Hamiltonian system even for long integration times, at least if the step size is small enough that the expansion converges. Although the approximate dynamics of a symplectic integrator will diverge if the step size is too large, the inherent stability of this approximation admits much larger step sizes, and hence reduced computation, than other discretizations of the dynamics.

This symmetric leapfrog integrator enjoys a global error quadratic in the step size and consequently it is classified as a second-order integrator. Higher-order integrators can just as easily be built up by applying each component operator multiple times in careful arrangements to cancel more and more error terms.

#### 4. EXPERIMENTS

To explore the performance of symplectic optimization numerically, we consider a relatively simple experiment. Let  $\mathcal{X}$  be a 50-dimensional Euclidean space equipped with the quadratic objective function

$$f(x) = \langle \Sigma^{-1}x, x \rangle,$$

where

$$\Sigma_{ij} = \rho^{|i-j|},$$

and  $\rho = 0.9$  to correlate the objective.

Let the auxiliary function,  $h$ , also be quadratic but without any interactions among the coordinates,

$$h(x) = \langle x, x \rangle.$$

The Bregman Hamiltonian becomes

$$H(x, r, t) = \frac{1}{2}e^{\alpha(t)-\gamma(t)} \langle p, p \rangle + e^{\alpha(t)+\beta(t)+\gamma(t)} f(x),$$

and the component vector fields of the extended dynamics take the form

$$\begin{aligned} \vec{H}_A &= \frac{d}{dt} \\ \vec{H}_{B1} &= 0 \\ \vec{H}_{B2} &= -\frac{1}{2} \left( \frac{\partial \alpha}{\partial t}(t) - \frac{\partial \gamma}{\partial t}(t) \right) e^{\alpha(t)-\gamma(t)} \langle p, p \rangle \frac{d}{d\mathcal{E}} \\ \vec{H}_{B3} &= e^{\alpha(t)-\gamma(t)} r \frac{d}{dx} \\ \vec{H}_{C1} &= -e^{\alpha(t)+\beta(t)+\gamma(t)} \nabla f(x) \frac{d}{dr} \\ \vec{X}_{C2} &= - \left( \frac{\partial \alpha}{\partial t}(t) + \frac{\partial \beta}{\partial t}(t) + \frac{\partial \gamma}{\partial t}(t) \right) e^{\alpha(t)+\beta(t)+\gamma(t)} f(x) \frac{d}{d\mathcal{E}}. \end{aligned}$$

In this case each of these vector fields are trivial and hence can be integrated exactly.

Finally we adopt the ideal scaling conditions for  $\alpha(t)$ ,  $\beta(t)$  and  $\gamma(t)$  discussed in Section 2, in which case the vector fields become

$$\begin{aligned}\vec{H}_A &= \frac{d}{dt} \\ \vec{H}_{B1} &= 0 \\ \vec{H}_{B2} &= \frac{1}{2} \frac{p(p+1)}{t^{p+2}} \langle r, r \rangle \frac{d}{d\mathcal{E}} \\ \vec{H}_{B3} &= \frac{p}{t^{p+1}} r \frac{d}{dx} \\ \vec{H}_{C1} &= -C p t^{2p-1} \nabla f(x) \frac{d}{dr} \\ \vec{H}_{C2} &= -C p (2p-1) t^{2p-2} f(x) \frac{d}{d\mathcal{E}}.\end{aligned}$$

Applying these component dynamics to the second-order leapfrog integrator introduced in Section 3.3 then gives the symmetric update sequence

$$\begin{aligned}t_{n+\frac{1}{2}} &= t_n + \epsilon \\ \mathcal{E}_{n+\frac{1}{2}} &= \mathcal{E}_n + \epsilon \left( \frac{1}{2} \frac{p(p+1)}{t^{p+2}} \langle r_n, r_n \rangle + -C p (2p-1) t^{2p-2} f(x_n) \right) \\ r_{n+\frac{1}{2}} &= r_n - \epsilon C p t^{2p-1} \nabla f(x_n) \\ x_{n+1} &= x_n + \epsilon \frac{p}{t^{p+1}} r_{n+\frac{1}{2}} \\ r_{n+1} &= r_{n+\frac{1}{2}} - \epsilon C p t^{2p-1} \nabla f(x_{n+1}) \\ \mathcal{E}_{n+1} &= \mathcal{E}_{n+\frac{1}{2}} + \epsilon \left( \frac{1}{2} \frac{p(p+1)}{t^{p+2}} \langle r_{n+1}, r_{n+1} \rangle + -C p (2p-1) t^{2p-2} f(x_{n+1}) \right) \\ t_{n+1} &= t_{n+\frac{1}{2}} + \epsilon.\end{aligned}$$

For comparison we implement the three-step dynamical Nesterov discretization derived in [Wibisono et al. \(2016\)](#). Given the ideal scaling conditions and the quadratic auxiliary function,  $h$ , this algorithm is given by

$$\begin{aligned}x_{n+1} &= \frac{p}{n+1} z_n + \left(1 - \frac{p}{n+1}\right) y_n \\ y_{n+1} &= x_{n+1} - \frac{p \epsilon^p}{2N} \nabla f(x_{n+1}) \\ z_{n+1} &= z_n - \epsilon^p C p (n+1)^{p-1} \nabla f(y_{n+1}).\end{aligned}$$

For both the extended Hamiltonian system and the Nesterov sequence we take  $\epsilon = 0.1$ ,  $p = 2$ ,  $C = 0.0625$ , and  $N = 2$ .

Results of this experiment are shown in Figure 1a. We see that the initial convergence rate obtained by symplectic integration and the three-step generalized Nesterov discretization are both roughly  $\mathcal{O}(t^{-2.95})$  for this problem. This should be no surprise, given that both approaches are stable discretizations of the same underlying Bregman dynamics. The number of iterations to arrive near the optimum is accordingly similar for the two algorithms for large values of the error criterion. It is smaller for the Nesterov discretization in the case of smaller error values. We return to this phenomenon—the increasing rate of the Nesterov discretization as it approaches the optimum—in the following section.

It is important to emphasize that the number of iterations is not the same as wall-clock time. Indeed, the leapfrog integrator that drives our implementation of symplectic optimization requires only a single gradient evaluation per iteration while the three-step generalized Nesterov discretization requires two to achieve stability. Consequently, as shown in Figure 1b, once we normalize for computational cost the symplectic integrator becomes twice as effective for large values of the error criterion. Whether this improvement persists in comparison to two-step Nesterov algorithms is an open question.

It is also important to emphasize that with the leapfrog integrator we are able to choose larger step sizes than with the three-step Nesterov discretization. This is due to the inherent stability of symplectic integration. In particular, as shown in Figure 1c, if we increase the step size to  $\epsilon = 0.25$  we see that the Nesterov discretization quickly diverges while the leapfrog integrator remains stable. The time to arrive near the optimum decreases uniformly for the leapfrog integrator for this larger value of the step size.

## 5. ACHIEVING EXPONENTIAL CONVERGENCE WITH A GRADIENT FLOW

As we have seen in Figure 1, the three-step generalized Nesterov discretization exhibits a unique behavior once it has become sufficiently close to the minimum. In that neighborhood the dynamical Nesterov discretization transitions into an exponential rate of convergence towards the minimum and soon surpasses the symplectic optimizer. Interestingly, we have found that this behavior does not persist for a quartic objective function,  $f(x) = \langle x, x \rangle^2$ , suggesting that it requires strong convexity of the neighborhood of the objective. Exponential convergence of the generalized Nesterov discretization in regions of strong convexity of the objective was considered in Wibisono et al. (2016).

Because this phase of exponential convergence does not appear in symplectic optimization it cannot be a feature of the Bregman dynamics themselves. Instead it must be a side effect of the heuristic discretization of the generalized Nesterov discretization, which introduced the auxiliary sequence,

$$y_n = \operatorname{argmin}_{y \in \mathcal{X}} \left[ f_{p-1}(y, ; x_n) + \frac{N}{\epsilon^p p} \|y - x_n\|^p \right],$$

or, for the conditions of Section 4,

$$y_{n+1} = x_{n+1} - \frac{p \epsilon^p}{2N} \nabla f(x_{n+1}).$$

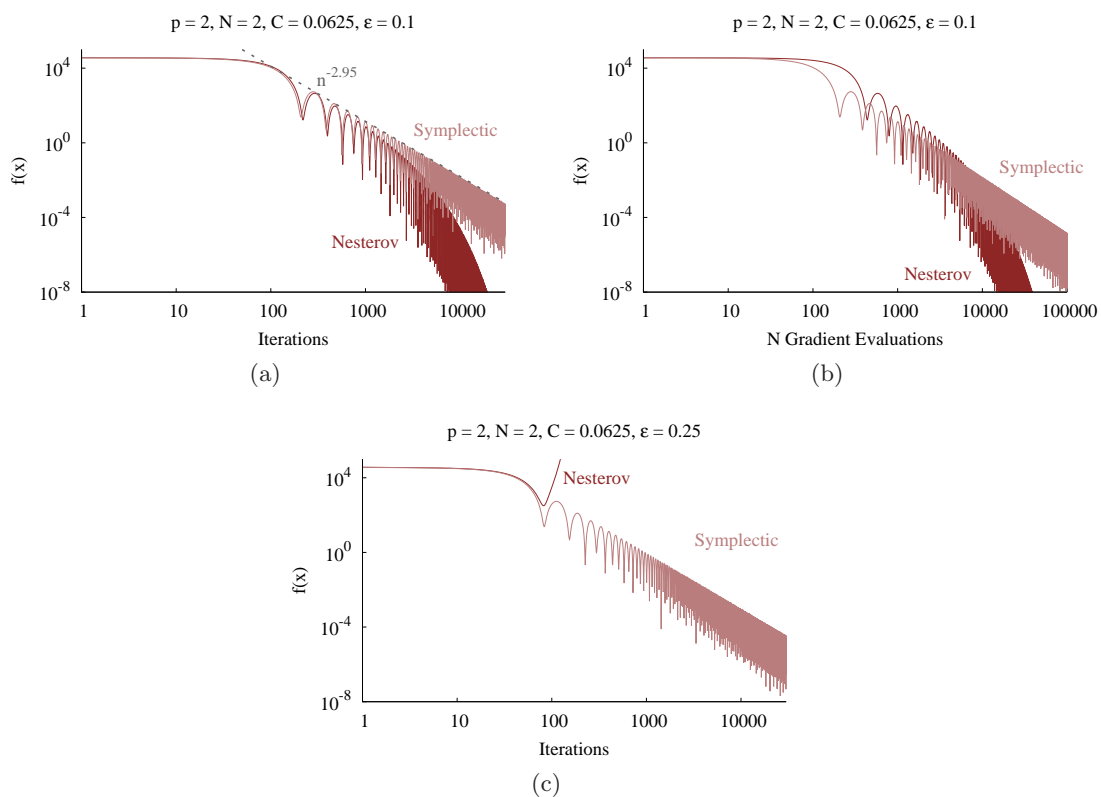


FIG 1. (a) When appropriately tuned, both symplectic optimization and the dynamic Nesterov discretization simulate the same latent Bregman dynamics and hence achieve similar convergence rates, here approximately  $\mathcal{O}(t^{-2.95})$ . (b) The symplectic optimization, however, requires only half of the computational effort of the three-step generalized Nesterov discretization. (c) Moreover, the inherent stability of the symplectic optimization admits larger discretization step sizes and even higher performance improvements.

For these conditions this sequence actually simulates a gradient flow on the configuration manifold, and consequently its addition interweaves the Bregman dynamics with a gradient flow. The exact nature of this interweaving, however, seems to ensure that the two evolutions characterize the dynamic Nesterov discretization in different regimes. Away from the minimum of the objective the Bregman dynamics dominate, rapidly pulling the system towards the minimum. Asymptotically, however, the dynamics dampen and eventually the gradient flow becomes dominant.

For sufficiently well-behaved objectives the emergence of the gradient flow allows admits the exponential convergence seen in the quadratic objective of Section 4. The gradient flow not only not only stabilizes the dynamic Nesterov discretization, it can also provide for even faster convergence near the minimum of the objective!

Although symplectic optimization doesn't need a gradient flow for stability, it could possibly benefit from the potentially exponential convergence it admits. Fortunately, incorporating a gradient flow into a symplectic integrator is straightforward—instead of trying to approximate the evolution operator  $\exp(\epsilon \vec{H}_\Xi)$  we instead try to approximate

$$\exp\left(\epsilon\left(\vec{H}_\Xi + \vec{X}_{\text{GF}}\right)\right),$$

where

$$\vec{X}_{\text{GF}} = -\frac{p \epsilon^p}{2N} \nabla f(x) \frac{d}{dx}$$

is the gradient field generating the gradient flow. Provided that we construct an appropriate symmetric splitting then the resulting integrator will enjoy the same global error as a symplectic integrator applied to the extended Hamiltonian system.

For the leapfrog integrator we constructed in Section 3.3 all we have to do is add  $\vec{X}_{\text{GF}}$  to the central operator, replacing  $\exp(\epsilon(\vec{H}_{B3}))$  with  $\exp(\epsilon(\vec{H}_{B3} + \vec{H}_{\text{GF}}))$ . Although this combined evolution operator is technically nonlinear and requires an implicit solution, here we will approximate the evolution with the explicit update,

$$x_{n+1} = \frac{p}{t^{p+1}} r_{n+\frac{1}{2}} - \frac{p \epsilon^p}{2N} \nabla f(x_{n-1}).$$

In the quadratic case where the dynamic Nesterov discretization exhibited exponential convergence, this modification of symplectic optimization exhibits the same advantageous behavior (Figure 2a). Moreover, the modified symplectic optimization maintains its superior stability, still allowing for larger step sizes and faster practical convergence (Figure 2b).

Still, while the global error scaling is preserved under the addition of the gradient flow and the modified Hamiltonian optimization works well empirically, we cannot always expect the same stability with the gradient flow. The problem is that gradient flow cannot be generated from a Hamiltonian and hence the modified discretized evolution cannot preserve the symmetries of the underlying dynamics. In practice we have to be careful to tune

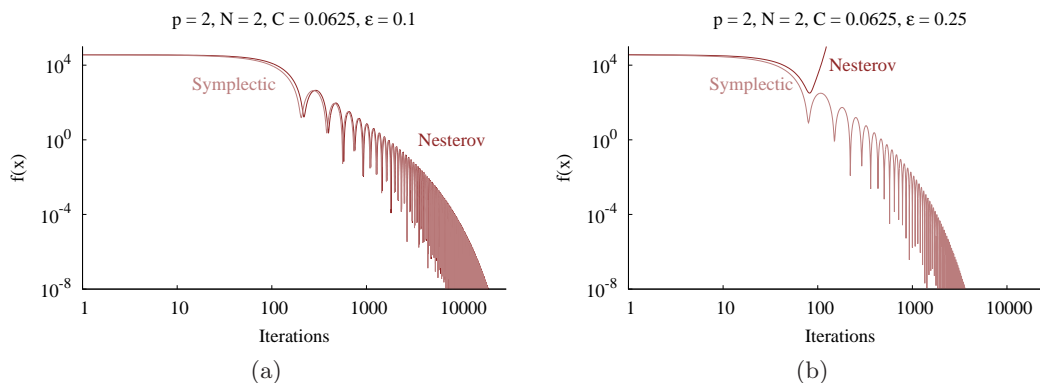


FIG 2. (a) By incorporating gradient flow into the leapfrog integration of the Bregman Hamiltonian dynamics we recover the same asymptotic exponential convergence near the minimum of the objective exhibited by the generalized Nesterov discretization. (b) These modified Hamiltonian dynamics remain stable even as we increase the step size, allowing for more efficient computation without compromising the advantageous asymptotic behavior.

the gradient flow so that its contributions, including any violations of the dynamical symmetries, are negligible until the Hamiltonian dynamics have converged close to the minimum of the objective.

## 6. DISCUSSION

Wibisono et al. (2016) introduced a dynamical system that converged to the minimum of a given objective function at the same rate as accelerated Nesterov methods. Moreover, by carefully discretizing the Lagrangian representation of these dynamics they were able to explicitly derive entire families of known accelerated Nesterov discretizations. Given the dynamical system itself, however, discretization is more systematically achieved by considering the Hamiltonian view of the system and appealing to symplectic integrators.

In particular, this systematic approach allows us to isolate the effects of the dynamics from other modifications, such as the gradient flow added to stabilize the original discretization of the Lagrangian representation of the dynamics. This separation then allows us to analyze the performance of the latent Bregman dynamics and that of any amendments independently.

This then positions us to study the general nature and optimality of the Bregman dynamics themselves. This study may not even be limited to Euclidean configuration spaces but perhaps also any manifold in a single unified setting. We discuss details of the systematically geometric construction of the Bregman dynamics and possible generalizations in Appendix A.

This systematic foundation may also allow us to formalize many of the empirical behaviors exhibited by Nesterov methods. For example, the folk wisdom is that Nesterov

methods do not perform particularly well when the objective is stochastic. This behavior, however, is not particularly surprising given the nature of symplectic integrators. As discussed in [Betancourt \(2015\)](#), the stochastic variations in the objective introduces a bias into symplectic integrators that corrupts their accuracy by pushing the numerical approximations away from the true dynamics. Intuitively the dynamical evolution moves so quickly that the variation in the stochastic objective doesn't have sufficient time to average out, unlike slower methods such as such as Robbins-Monro that do work well with stochastic objectives. On the other hand, the time dependence of the Bregman dynamics may provide a way of compensating for this bias. Only with a formal understanding of the Bregman dynamics afforded by this new perspective will be able to identify the necessary structure.

Unfortunately, the introduction of Hamiltonian symplectic integrators also complicates the formal analysis of Hamiltonian optimization itself. For example, the accuracy of leapfrog integrators comes from cancellations in their symmetric updates, but any individual update can have large error. Hence we cannot expect to be able to bound convergence term-by-term. Indeed the stability of symplectic integrators is a global property—the discretized dynamics oscillate around the true dynamics and discrete updates will in general deviate away from the exact dynamics before finally returning. To understand the convergence of the discretized dynamics we instead have to take non-local and topological considerations into account, as is done in *backwards error analysis* ([McLachlan et al., 2004](#); [Leimkuhler & Reich, 2004](#); [Hairer et al., 2006](#)).

Ultimately, however, the direct window into Bregman dynamics provided by their Hamiltonian representation and corresponding symplectic integration enables not only a better understanding of existing accelerated Nesterov methods but also a principled way of developing new implementations and generalizations.

## 7. ACKNOWLEDGEMENTS

We thank Sam Power for helpful comments.

## REFERENCES

- Betancourt, M. The fundamental incompatibility of scalable Hamiltonian Monte Carlo and naive data subsampling. In *Proceedings of the 32nd International Conference on Machine Learning*, pp. 533–540, Lille, France, 2015.
- de León, M. and Rodrigues, P. R. *Methods of Differential Geometry in Analytical Mechanics*, volume 158 of *North-Holland Mathematics Studies*. North-Holland Publishing Co., Amsterdam, 1989.
- Hairer, E., Lubich, C., and Wanner, G. *Geometric Numerical Integration: Structure-Preserving Algorithms for Ordinary Differential Equations*. Springer, New York, 2006.
- Jin, C., Netrapalli, P., and Jordan, M. I. Accelerated gradient descent escapes saddle points faster than gradient descent. *Arxiv preprint Arxiv:1711.10456*, 2017.
- José, J. V. and Saletan, E. J. *Classical Dynamics: A Contemporary Approach*. Cambridge University Press, New York, 1998.
- Krichene, W., Bayen, A., and Bartlett, P. Accelerated mirror descent in continuous and discrete time. In *Advances in Neural Information Processing Systems (NIPS) 27*, 2015.
- Lee, J. M. *Introduction to Smooth Manifolds*. Springer, 2013.

- Leimkuhler, B. and Reich, S. *Simulating Hamiltonian Dynamics*. Cambridge University Press, New York, 2004.
- McLachlan, R. I., Perlmutter, M., and Quispel, G. R. W. On the nonlinear stability of symplectic integrators. *BIT Numerical Mathematics*, 44(1):99–117, 2004.
- Nemirovskii, A. and Yudin, D. *Problem Complexity and Method Efficiency in Optimization*. John Wiley, New York, 1983.
- Nesterov, Y. A method of solving a convex programming problem with convergence rate  $o(1/k^2)$ . *Soviet Mathematics Doklady*, 27:372–376, 1983.
- Su, W., Boyd, S., and Candes, E. J. A differential equation for modeling Nesterov’s accelerated gradient method: theory and insights. *Journal of Machine Learning Research*, 17(153):1–43, 2016.
- Wibisono, A., Wilson, A. C., and Jordan, M. I. A variational perspective on accelerated methods in optimization. *Proceedings of the National Academy of Sciences*, 113(47):E7351–E7358, 2016.

## A. GEOMETRIC CONSTRUCTION OF THE BREGMAN DYNAMICS

Symplectic integrators are not only straightforward to implement and extremely powerful in practice, they are applicable to Hamiltonian dynamics defined over any manifold. This motivates the consideration as to whether or not the Bregman dynamics themselves could be generalized beyond a Euclidean configuration space and onto a more general manifold.

In the main paper we heavily utilized the canonical coordinates of a Euclidean manifold in the construction of the Bregman Hamiltonian, paying relatively little attention to the difference between points in the configuration space,  $\mathcal{X}$ , elements of its local tangent and cotangent spaces, and even its tangent and cotangent bundles. Here we take a systematic geometric perspective in an attempt to highlight some possible paths towards generalizing symplectic optimization.

As in the main text we begin with a Euclidean manifold,  $\mathcal{X}$ , equipped with the auxiliary function,  $h$ . We have to be careful, however, as there are two ways of defining  $h$  that lead to the same Bregman dynamics, and these two approaches may prove to motivate different generalizations.

### A.1 Notation

We will largely follow the notation of Lee (2013), although contraction of a vector field with a one-form will be denoted with  $\lrcorner$ .

For any manifold,  $\mathcal{X}$ , we can construct an associated tangent bundle,  $\pi : T\mathcal{X} \rightarrow \mathcal{X}$ . Correspondingly, any point in the tangent bundle,  $w \in T\mathcal{X}$  identifies both a point in the base manifold,  $\pi(w) \in \mathcal{X}$  and a vector in the local tangent space,  $v(x) \in T_{\pi(w)}\mathcal{X}$ .

We can also construct the dual cotangent bundle,  $\varpi : T^*\mathcal{X} \rightarrow \mathcal{X}$ . Any point in the cotangent bundle,  $z \in T^*\mathcal{X}$ , identifies both a point in the base manifold,  $\varpi(z) \in \mathcal{X}$  and a covector in the local cotangent space,  $p(x) \in T_{\varpi(z)}^*\mathcal{X}$ .

### A.2 Defining the auxiliary function on the configuration space

The first way to proceed is to define the auxiliary function directly on the configuration manifold,

$$h : \mathcal{X} \rightarrow \mathbb{R}.$$



The corresponding Bregman Divergence is a bit ungainly to write geometrically, requiring liberal use of the identification of  $\mathcal{X}$  with any tangent space, but the kinetic energy defining the dynamical system has much cleaner form.

Given any point in the tangent bundle,  $w \in T\mathcal{X}$ , we can construct a new point in the base manifold by translating along the Euclidean connection for some time,  $T$ ,

$$\pi(w) + T v(w) = x' \in \mathcal{X}.$$

A natural question is how the auxiliary function,  $h$ , changes under this parallel transport,

$$\Delta h(w, T) = h(\pi(w) + T v(w)) - h(\pi(w)),$$

relative to how it would change differentially,

$$\delta h(w, T) = (\mathrm{d}h \lrcorner T v(w))(\pi(w)) = T (\mathrm{d}h \lrcorner v(w))(\pi(w)).$$

Setting  $T = \exp(-\alpha(t))$  we can then define the kinetic energy as a function on the time-dependent tangent bundle,

$$K : T\mathcal{X} \times \mathbb{R} \rightarrow \mathbb{R},$$

where

$$\begin{aligned} K(w, t) &= \Delta h(w, e^{-\alpha(t)}) - \delta h(w, e^{-\alpha(t)}) \\ &= h(\pi(w) + e^{-\alpha(t)} v(w)) - h(\pi(w)) \\ &\quad - e^{-\alpha(t)} (\mathrm{d}h \lrcorner v(w))(\pi(w)). \end{aligned}$$

This kinetic energy allows us to construct the time-dependent Bregman Lagrangian,

$$\begin{aligned} L : T\mathcal{X} \times \mathbb{R} &\rightarrow \mathbb{R} \\ (w, t) &\mapsto e^{\alpha(t)+\gamma(t)} (K(w, t) - U(\pi(w))), \end{aligned}$$

where, as before,

$$\begin{aligned} U : \mathcal{X} \times \mathbb{R} &\rightarrow \mathbb{R} \\ (x, t) &\mapsto e^{\beta(t)} f(x). \end{aligned}$$

Taking the natural coordinates for  $\mathcal{X}$  this reduces to exactly the Bregman Lagrangian discussed in the main text, hence they are equivalent functions.

We can now build the Bregman Hamiltonian as the Legendre dual of  $L$  on  $T^*\mathcal{X} \times \mathbb{R}$ . As before we define the components of the conjugate momenta as

$$p^i = \frac{\partial L}{\partial v_i}(x, v, t),$$

with the Hamiltonian given by the Legendre transform

$$\begin{aligned} H : T^* \mathcal{X} \times \mathbb{R} &\rightarrow \mathbb{R} \\ (x, p, t) &\mapsto p(v(x, p, t)) - L(x, v(x, p, t), t). \end{aligned}$$

In the main text we were able to solve for the velocities as a function of the momenta and cleanly substitute then into the Bregman Hamiltonian with the Legendre transform of the auxiliary function,  $h$ . We have to be careful here, however, because as  $h$  is not defined on the tangent bundle it does not admit the same Legendre transform as the Bregman Lagrangian.

Instead we must exploit the Euclidean structure of  $\mathcal{X}$ . Because  $\mathcal{X}$  is a vector space with dual  $\mathcal{Y}$ , we can perform a Legendre transform from  $\mathcal{X}$  to  $\mathcal{Y}$  to define

$$h^* : \mathcal{Y} \rightarrow \mathbb{R}.$$

On a Euclidean manifold every element of  $\mathcal{Y}$  is identified with an element of any cotangent space, so this conjugate also defines a function

$$\eta^* : T^* \mathcal{X} \rightarrow \mathbb{R}.$$

The gradient of  $h^*$ , however, is identified with a vector field over  $\mathcal{X}$ . In particular, the gradient of  $h^*$  is equivalent to the Euler vector field,

$$E = x^i \frac{\partial}{\partial x^i}.$$

Utilizing these two objects we can then write the Bregman Hamiltonian geometrically as

$$\begin{aligned} H(z, t) = & e^{\alpha(t)+\gamma(t)} \left( g^*(e^{-\gamma(t)} p(z) + dh(\varpi(z))) \right. \\ & - g^*(dh(\varpi(z))) - e^{-\gamma(t)} (p(z) \lrcorner E)(\pi(z)) \\ & \left. + e^\beta f(\pi(z)) \right). \end{aligned}$$

Once again, in Euclidean coordinates this reduces to what was presented in the main text.

### A.3 Defining the auxiliary function on the tangent bundle

The other way to proceed is to define the auxiliary function on the tangent bundle,

$$h : T\mathcal{X} \rightarrow \mathbb{R}.$$

The corresponding Bregman Divergence is still a bit ungainly to write geometrically, but the kinetic energy also exhibits a cleaner form.

In this case we exploit the identification of every point  $x \in \mathcal{X}$  with a vector  $u(x) \in T_x\mathcal{X}$ . A point in the tangent bundle,  $w$ , then defines two vectors which we can add,

$$u(\pi(w)) + T v(w) = u' \in T_x\mathcal{X}.$$

As before we can write

$$\Delta h(w, T) = h(u(\pi(w)) + T v(w)) - h(u(\pi(w))),$$

only now  $\Delta h$  is a function on the tangent bundle.

The differential change in  $h$ , however, is a bit more subtle. Because  $h$  is now defined as a function on the tangent bundle, its differential,  $dh$  is a section of  $T^*T\mathcal{X}$ . In order to contract against the differential we need a section of  $TT\mathcal{X}$ . Fortunately we can use the rigid Euclidean connection of  $\mathcal{X}$  to define a unique horizontal lift from the vector field  $v(w)$  over  $\mathcal{X}$  to a vector field  $\tilde{v}(w)$  on  $T\mathcal{X}$ . Using this lift we can write

$$\delta h(w, T) = (dh \lrcorner T \tilde{v}(w))(w) = T (dh \lrcorner \tilde{v}(w))(w),$$

where again  $\delta h$  is now a function on the tangent bundle.

Putting these together and taking  $T = \exp(-\alpha(t))$  we get the time-dependent kinetic energy

$$K : T\mathcal{X} \times \mathbb{R} \rightarrow \mathbb{R}$$

where

$$\begin{aligned} K(w, t) &= \Delta h(w, \exp(-\alpha(t))) - \delta h(w, \exp(-\alpha(t))) \\ &= h(u(\pi(w)) + e^{-\alpha(t)} v(w)) - h(u(\pi(w))) \\ &\quad - e^{-\alpha(t)} (dh \lrcorner \tilde{v}(w))(w). \end{aligned}$$

Superficially this looks the same as the kinetic energy derived in Section A.2 but each term has a subtly different geometric interpretation. Still, once we impose Euclidean coordinates we see that the function, as well as the corresponding Lagrangian, is equivalent.

Continuing to the Hamiltonian, the construction of a Legendre conjugate for  $h$  is simplified as we simply need to use the same Legendre transform from the tangent to the cotangent bundle that we use to map the Lagrangian into a Hamiltonian. This gives

$$h^* : T^*\mathcal{X} \rightarrow \mathbb{R},$$

The gradient  $dh^*$  is now a section of  $T^*T^*\mathcal{X}$  which contracts against sections of  $TT^*\mathcal{X}$ . In order to construct such a section we can utilize rigid Euclidean structure of  $\mathcal{X}$  once again. On the cotangent bundle we have the canonical one-form  $\theta$  which we can raise to a vector field with a musical isomorphism,  $\theta^\sharp(z)$ . We can then define a unique horizontal lift of this vector field to the vector field  $\tilde{\theta}^\sharp(z)$  on  $TT^*\mathcal{X}$  using the Euclidean connection.

Similarly, for any element,  $z \in T^*\mathcal{X}$ , of the cotangent bundle we can identify an element of the tangent bundle by applying a musical isomorphism to the canonical one-form to give  $\theta_b(z) \in T\mathcal{X}$ .

With these objects we can now write the Bregman Hamiltonian in the alternative form

$$\begin{aligned} H(z, t) = & e^{\alpha(t)+\gamma(t)} \left( h^*(e^{-\gamma(t)}p(z) + dh(\theta_b(z))) \right. \\ & - h^*(dh(\theta_b(z))) - e^{-\gamma(t)}(dh^* \lrcorner \tilde{\theta}^\sharp)(z) \\ & \left. + e^\beta f(\pi(z)) \right). \end{aligned}$$

Once again, we have a Hamiltonian whose terms have a subtly different geometric interpretation but yield an equivalent Hamiltonian function on the cotangent bundle.

As we're utilizing the same Euclidean structure of  $\mathcal{X}$  these two approaches shouldn't give contradictory answers, and we see here that they don't. The two paths, however, may illuminate different strategies for generalizing the construction of a Bregman dynamical system beyond the Euclidean case. For example, the second approach seems particularly appropriate for Riemannian manifolds equipped with an appropriate connection.

Regardless of how we might generalize or modify the Bregman dynamics, provided we maintain a geometric construction the discretization will proceed as smoothly as in the Euclidean case thanks to the geometric universality of symplectic integrators.