

# Lagrangian and Hamiltonian Dynamics on the Simplex

Goffredo Chirco, Luigi Malagò, Giovanni Pistone

{chirco,malago}@rist.ro, giovanni.pistone@carloalberto.org

## Abstract

The statistical bundle is the set of couples  $(q, w)$  of a positive probability density  $q$  and a random variable  $w$  such that  $\mathbb{E}_q[w] = 0$ . On a finite state space, we assume  $q$  to be a probability density with respect to the uniform probability and express it in the affine atlases of exponential charts. Velocity and acceleration of a one-dimensional statistical model are computed using the canonical dual pair of parallel transports. We define Lagrangian and Hamiltonian mechanics on the bundle and we provide explicit examples of a time-independent and time-dependent Lagrangian functions, leading to diverse accelerated natural gradient dynamics. Within this formalism we can reproduce Nesterov's flow for convex constrained optimization problems on the statistical bundle.

## Statistical Bundle

We consider a finite sample space  $\Omega$  with cardinality  $N$ . Let  $\Delta(\Omega)$  be the probability simplex, and  $\Delta^\circ(\Omega)$  its interior. We denote with  $\mu$  the uniform probability function  $1/N$ .

The **maximal exponential family**  $\mathcal{E}(\mu)$  is the set of densities which can be written as  $p \propto e^f$ , where  $f$  is defined up to a constant. Given a reference density  $p \in \mathcal{E}(\mu)$ , we have

$$q(x) = \exp(v(x) + H(v)) \cdot p(x),$$

with  $\mathbb{E}[v(x)] = 0$ ,  $H(v) = -\log \mathbb{E}_q[e^v] = D(p \| q)$  and

$$v = \log \frac{q}{p} - \mathbb{E}_q \left[ \log \frac{q}{p} \right].$$

The **exponential statistical bundle** with base  $\Omega$  is defined as

$$S\mathcal{E}(\mu) = \{(q, v) \mid q \in \mathcal{E}(\mu), \mathbb{E}_q[v] = 0\},$$

we denote with  ${}^*S_q\mathcal{E}(\mu)$  the **dual statistical bundle**. For finite  $\Omega$ ,  $S_q\mathcal{E}(\mu)$  and  ${}^*S_q\mathcal{E}(\mu)$  coincide.

## Affine Geometries

A duality mapping between the statistical bundle and its dual the can be defined at the fiber at  $q$  by

$${}^*S_q\mathcal{E}(\mu) \times S_q\mathcal{E}(\mu) \ni (\eta, v) \mapsto \langle \eta, v \rangle_q = \mathbb{E}_q[\eta v].$$

Two different affine geometries can be define for  $S_q\mathcal{E}(\mu)$  and  ${}^*S_q\mathcal{E}(\mu)$ , by defining two different transports for each  $p, q \in \mathcal{E}(\mu)$ , i.e.,

- Exponential transport**  $e\mathbb{U}_p^q: S_p\mathcal{E}(\mu) \rightarrow S_q\mathcal{E}(\mu)$ ,  $e\mathbb{U}_p^q v = v - \mathbb{E}_q[v]$ ,
- Mixture transport**  $m\mathbb{U}_p^q: {}^*S_p\mathcal{E}(\mu) \rightarrow {}^*S_q\mathcal{E}(\mu)$ ,  $m\mathbb{U}_p^q \eta = \frac{q}{p} \eta$ .

## Quadratic Lagrangian

Let  $m$  be the inertial mass

$$L(q, w) = \frac{m}{2} \mathbb{E}_q[w^2] = \frac{m}{2} \langle w, w \rangle_q, \quad m \geq 0, (q, w) \in S\mathcal{E}(\mu).$$

We can obtain an expression in the chart centered in  $p$  for the Lagrangian

$$L_p(u, v) = \frac{m}{2} \left\langle e\mathbb{U}_p^{e_p(u)} v, e\mathbb{U}_p^{e_p(u)} v \right\rangle_{e_p(u)} = \frac{m}{2} d^2 K_p(u)[v, v],$$

where  $q = e_p(u)$  and  $w = e\mathbb{U}_p^q v$

By computing the total derivative in the chart of  $L$ ,

$$dL_p(u, v)[h, k] = \frac{m}{2} \left\langle w^2 - \mathbb{E}_q[w^2], e\mathbb{U}_p^q u \right\rangle_q + m \langle w, e\mathbb{U}_p^q k \rangle_q.$$

we can obtain the Euler-Lagrange equation

$$\frac{D}{dt} \dot{q}(t) = \frac{1}{2} \left( \dot{q}(t)^2 - \mathbb{E}_{q(t)}[\dot{q}(t)^2] \right),$$

which can be expressed as a system of  $N$  second-order ODEs

$$\ddot{q}_j(t) = \frac{\dot{q}_j(t)^2}{2q_j(t)} - \frac{q_j(t)}{2N} \sum_{i=1}^N \frac{\dot{q}_i(t)^2}{q_i(t)^2}, \quad j = 1, \dots, N.$$

Consider the case of a Lagrangian function given by the difference of the quadratic form and a potential on the bundle,

$$L(q, w) = \frac{m}{2} \langle w, w \rangle_q - \kappa \mathbb{E}_q[\log q],$$

with the negative entropy  $f(q) = -\mathcal{H}(q)$  playing the role of the convex potential well.

The Euler-Lagrange equation can be derived as

$$m \frac{D}{dt} \dot{q} = \frac{m}{2} \left( \dot{q}(t)^2 - \mathbb{E}_{q(t)}[\dot{q}(t)^2] \right) + \kappa \text{grad } \mathcal{H}(q).$$

Let  $A(q, v) = v^2/2 + \frac{\kappa}{m} \log(q)$  and  $B(q, v) = v^2/2 - \frac{\kappa}{m} \log(q)$ , the associated system of first-order ODEs is

$$\begin{cases} \frac{d}{dt} q(x; t) = q(x; t) v(x; t) \\ \frac{d}{dt} v(x; t) = -A(q(x; t), v(x; t)) - \frac{1}{N} \sum_y q(y; t) B(q(y; t), v(y; t)) \end{cases}$$

for  $x \in \Omega$ .

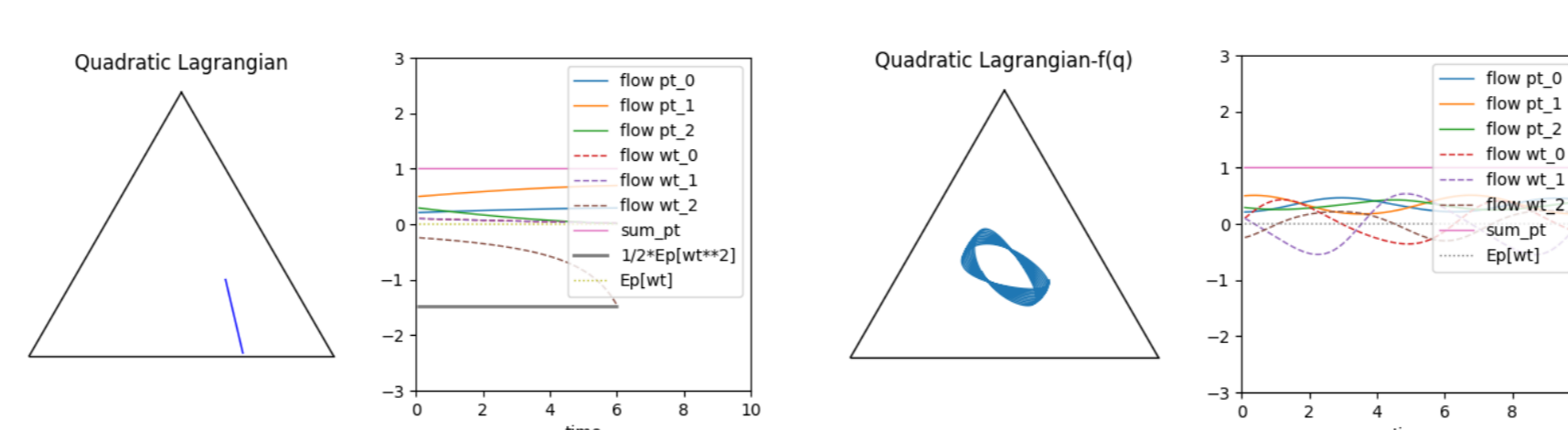


Figure 1: (Left) Free particle quadratic potential; (Right) Motion in Potential.

## Kullback-Leibler Lagrangian

A divergence is a smooth mapping  $D: \mathcal{E}(\mu) \times \mathcal{E}(\mu) \rightarrow \mathbb{R}$ , such that for all  $p, q \in \mathcal{E}(\mu)$  it holds  $D(p, q) \geq 0$  and  $D(p, q) = 0$  if, and only if,  $p = q$ .

Every divergence can be associated to a Lagrangian by the canonical mapping

$$\mathcal{E}(\mu)^2 \ni (q, r) \mapsto (q, s_q(r)) = (q, w) \in S\mathcal{E}(\mu),$$

with  $q = e^{v-K_p(v)} \cdot p$ , that is,  $v = s_p(q)$ .

We have an equivalence of a couple of a point and a vector and a couple of points. Every divergence  $D$  is mapped into a *divergence Lagrangian*, and conversely,

$$L(q, w) = D(q, e_q(w)), \quad D(q, r) = L(q, s_q(r)).$$

We focus on the case of the Kullback-Leibler divergence (KL), which lies at the intersection of the family of Csiszár's  $f$ -divergences and Bregman divergences (Amari, 2016).

Up to second-order approximation the KL provides a *locally quadratic measure*, motivating its interpretation as a local, non-symmetric generalization of the kinetic energy of classical mechanics.

The Lagrangian

$$D(q, r) = D(q \| r) = \mathbb{E}_q \left[ \log \frac{q}{r} \right],$$

can be written in chart at  $q$  as

$$D(q \| e_q(w)) = \mathbb{E}_q \left[ \log \frac{q}{e_q(w)} \right] = \mathbb{E}_q[-w + K_q(w)] = K_q(w).$$

The expression of the divergence Lagrangian in chart at  $p$  is

$$L_p(u, v) = L(e_p(u), e\mathbb{U}_p^{e_p(u)} v) = D(e_p(u), e_{e_p(u)}(e\mathbb{U}_p^{e_p(u)} v)) = D(e_p(u), e_p(u + v)).$$

The Euler-Lagrange equation is obtained by plugging in  $w(t) = \dot{q}(t)$ ,

$$\frac{D}{dt} \left( e^{\dot{q}(t)-K_{q(t)}(\dot{q}(t))} - 1 \right) = e^{\dot{q}(t)-K_{q(t)}(\dot{q}(t))} - 1 - \dot{q}(t),$$

which takes the form of a second-order equation

$$\left( e^{\dot{q}(t)-K_{q(t)}(\dot{q}(t))} \right) \left( \dot{q}(t) + \ddot{q}(t) - \mathbb{E}_{e_{q(t)}(\dot{q}(t))}[\dot{q}(t) + \ddot{q}(t)] \right) = e^{\dot{q}(t)-K_{q(t)}(\dot{q}(t))} - 1.$$

By using  $\dot{q}(t) = v(t)$  we have

$$\frac{d}{dt} v(t) = \ddot{q}(t) - \mathbb{E}_{q(t)}[v(t)^2] = -v(t) + \frac{e^{\dot{q}(t)-K_{q(t)}(\dot{q}(t))} - 1}{e^{\dot{q}(t)-K_{q(t)}(\dot{q}(t))}} + \mathbb{E}_{q(t)} \left[ \frac{e^{\dot{q}(t)-K_{q(t)}(\dot{q}(t))} - 1}{e^{\dot{q}(t)-K_{q(t)}(\dot{q}(t))}} \right] - \mathbb{E}_{q(t)}[v(t)^2]$$

The strong convexity of the KL generating function ensures the existence of an invertible Legendre transform, naturally allowing for a Hamiltonian formulation.

Using the equation for  $\text{grad}_e K_q(w)$  and its inverse the Legendre transform of  $w \mapsto K_q(w)$  is

$$\begin{aligned} H_q(\eta) &= \langle \eta, \log(1 + \eta) - \mathbb{E}_q[\log(1 + \eta)] \rangle_q + \\ &\quad - K_q \left( \log(1 + \eta) - \mathbb{E}_q[\log(1 + \eta)] \right) \\ &= \mathbb{E}_q[\eta \log(1 + \eta)] - \mathbb{E}_q[\log(1 + \eta)] = \mathbb{E}_q[(1 + \eta) \log(1 + \eta)]. \end{aligned}$$

In the chart at  $p$ ,  $q = e_p(u) = e^{u-K_p(u)}$ ,  $p, \eta = m\mathbb{U}_p^{e_p(u)} \zeta = e^{-u+K_p(u)} \zeta$ , so that

$$H_p(u, \zeta) = \mathbb{E}_{e_p(u)} \left[ (1 + m\mathbb{U}_p^{e_p(u)} \zeta) \log(1 + m\mathbb{U}_p^{e_p(u)} \zeta) \right] = \mathbb{E}_p \left[ (e^{u-K_p(u)} + \zeta) \log(1 + e^{-u+K_p(u)} \zeta) \right].$$

By taking the derivative wrt  $u$ , and going back to the original variables, the Hamilton equations are

$$\begin{cases} \frac{D}{dt} \eta(t) = \eta(t) - \log(1 + \eta(t)) + \mathbb{E}_{q(t)}[\log(1 + \eta(t))] \\ \dot{q}(t) = \log(1 + \eta(t)) - \mathbb{E}_{q(t)}[\log(1 + \eta(t))] \end{cases}$$

The solution curve and its derivatives can be expressed in the global space in which the dual bundle is embedded by

$$\frac{D}{dt} \eta(t) = \frac{\dot{q}(t)}{q(t)} \eta(t) + \dot{\eta}(t), \quad \dot{q}(t) = \frac{\dot{q}(t)}{q(t)},$$

so that the resulting system of ODEs becomes

$$\begin{cases} \dot{\eta}(x; t) = \eta(x; t) - (1 + \eta(x; t)) (\log(1 + \eta(x; t)) - \frac{1}{N} \sum_y q(y; t) \log(1 + \eta(y; t))) \\ \dot{q}(x; t) = q(x; t) \left( \log(1 + \eta(x; t)) - \frac{1}{N} \sum_y q(y; t) \log(1 + \eta(y; t)) \right) \end{cases}$$

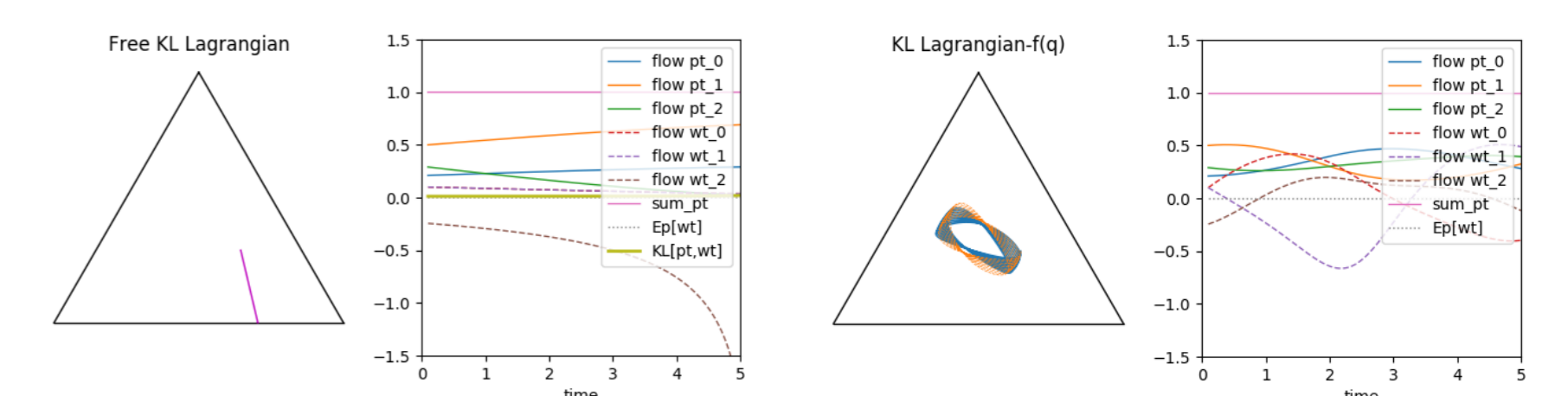


Figure 2: (Left) Free particle Kullback-Leibler divergence potential; (Right) Kullback-Leibler divergence motion in Potential.

## Time Dependent KL Lagrangian

We can introduce an explicit time dependence in the Lagrangian.

This choice is motivated by the role time in generating a *dissipative* accelerated dynamics, which is of interest in optimization.

In the exponential map, we consider a time-dependent scaling of the shift vector, such that  $\chi = e_q(e^{-\alpha t} w)$  and  $s_p(\chi) = u + e^{-\alpha t} v \in S_p\mathcal{E}(\mu)$ , with  $\alpha: I \rightarrow \mathbb{R}$  smooth,  $I \subset \mathbb{R}$  open time interval. With this choice the KL Lagrangian reads

$$D: I \times S\mathcal{E}(\mu) \ni (q, w, t) \mapsto D(q \| e_q(e^{-\alpha t} w)) \in \mathbb{R}.$$

In presence of explicit time-dependence, desirable closure under time-dilation can be achieved by an overall scaling of the divergence by a factor  $e^{\alpha t}$ , such that the new Lagrangian

$$L(q, w, t) = e^{\alpha t} D(q \| e_q(e^{-\alpha t} w)),$$

leads to fully time-reparametrization invariant action.

We can derive the Euler-Lagrange equation in presence of the time-scaling for  $v(t) = \dot{u}(t)$ , we get

$$\begin{aligned} d^2 K_p(u(t) + e^{-\alpha t} \dot{u}(t)) &= (e^{\alpha t} - \dot{\alpha} t) \dot{u}(t) + \ddot{u}(t), h = \\ &= e^{2\alpha t} (dK_p(u(t) + e^{-\alpha t} \dot{u}(t))[h] - dK_p(u(t))[h]), \end{aligned}$$

We can then transport the equation back on the statistical bundle to get

$$\frac{e_q(e^{-\alpha t} \dot{q})}{q} \left( (e^{\alpha t} - \dot{\alpha} t) \dot{q}(t) + \ddot{q}(t) - \mathbb{E}_{e_p(u+e^{-\alpha t} v)}[(e^{\alpha t} - \dot{\alpha} t) \dot{q}(t) + \ddot{q}(t)] \right) = e^{2\alpha t} \left( \frac{e_q(e^{-\alpha t} \dot{q})}{q} - 1 \right),$$

with respect to the equation derived for the cumulant Lagrangian, the time-dependent scaling leads to an extra *damping* contribution in the velocity, which redefines the coefficient of  $\dot{q}$ .

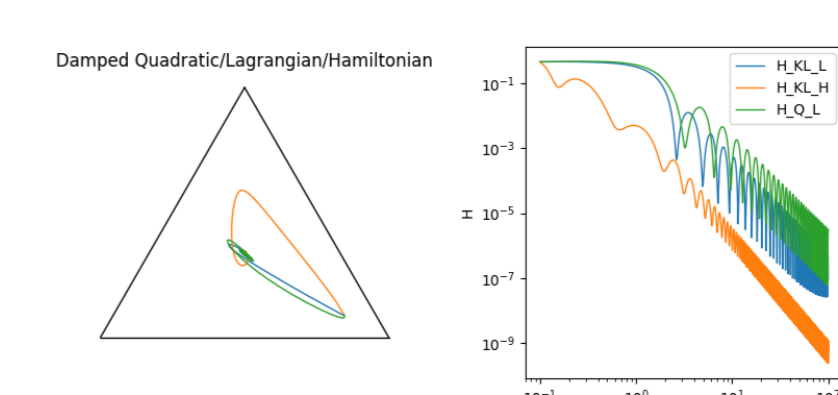


Figure 3: Comparison of different damped systems.