



CONVEXITY OF KULLBACK-LEIBLER FUNCTIONAL AND THE DEFINITION OF HESSIAN FOR REAL MAPS DEFINED ON FINITE MARKOV SPACES

Alireza BAHRAINI

Sharif University of Technology, Dept. of Mathematical Sciences, P.O.Box 11155-9415, Tehran, Iran.

E-mail: bahraini@sharif.edu

Abstract. In order to study the convexity of Kullback-Leibler functional in the presence of discrete variables, one needs to study the convexity of a linear functional defined on some type of Wasserstein space associated with these discrete spaces. Inspired by an observation in continuous case and by employing the Wasserstein spaces associated with Markov chains as introduced by Erbar-Maas [8], we define the Hessian of a real map defined on a Hypercube as the underlying Markov structure. The convexity of linear functionals in the discrete case can be computed through the Hessian matrices that can be derived through this relationship.

Keywords: finite Markov Spaces, Erbar-Maas Wasserstein space, geodesic convexity.

Mathematics Subject Classification (MSC2020): 39B62, 49Q20.

1. INTRODUCTION

Optimization in the space of probability measures is of special interest both in mathematics, statistics and machine learning. Many classical PDEs can be identified as minimization problems of specific functionals defined over Wasserstein type measure spaces [2, 4, 12]. In machine learning and statistics numerous problems including variational Bayesian inference, maximum likelihood estimations, generative adversarial networks, are stated as minimization over the space of probability distributions [5, 6].

The gradient flow of these probability functionals and their convergence rates towards critical points are governed by the geodesic convexity computed in terms of the geometry induced by Wasserstein metrics. This has led to novel numerical methods to solve PDEs of Fokker-Planck types, as well as new connections between the geometry of the finite dimensional underlying space and the associated Wasserstein probability space. A typical example consists of entropy functional and its role in developing a synthetic notion of Ricci curvature's lower bound for continuous and discrete spaces [9, 11].

Our goal in this note is to introduce a notion of convexity for maps defined on finite Markov spaces by applying similar ideas describing characteristics of an object of finite dimensional nature in terms of the associated Wasserstein space. In fact in the Kullback-Leibler (KL) functional $D_{KL}(p||q) = \int p \log p - \int p \log q$ we have two terms the entropy functional $\int p \log p$ and the simpler term $\int p \log q$. As mentioned above the convexity coefficient of the entropy term has been already discussed. We aim at providing some tools which enable us to investigate the geodesic convexity of the second term in the discrete case. KL functional is widely used in many statistical problems in particular in Mean Field Variational Inference (MFVI). The study of its critical points and convexity coefficients is of essential importance for exploring the rate of convergence of the corresponding mean field equations algorithms towards their fixed point. To this end we first observe in

section 2 that the convexity of a map $f : M \rightarrow \mathbb{R}$ defined over a Riemannian manifold (M, g) is equivalent to the geodesic convexity of the functional $\mathcal{L}^f : \mathcal{P}_2(\mathcal{M}) \rightarrow \mathbb{R}$ defined on the L^2 -Wasserstein space $\mathcal{P}_2(\mathcal{M})$ by

$$\mathcal{L}^f : \mathcal{P}_2(\mathcal{M}) \longrightarrow \mathbb{R}, \quad \mathcal{L}^f(\rho) := \int_{\mathcal{M}} f d\rho,$$

In the present note we introduce the Hessian of a real map defined on a hypercube. It turns out that one of the natural graphical structures for which a rich class of convex maps exists consists of the hypercube. This will lead to the uniqueness of solutions to mean field equations of the Ising model. By applying this observation one can compute new partition functions which had remained intractable. The main tool for us is the theory of Wasserstein spaces developed by Maas-Erbar for finite Markov spaces [9].

2. CONVEXITY OF REAL MAPS ON FINITE MARKOV SPACES THROUGH WASSERSTEIN SPACE

A subset $S \subset M$ of a Riemannian manifold (M, g) is geodesically convex if for all $x, y \in S$ there exists a geodesic $\gamma : [0, 1] \rightarrow M$, such that $\gamma(0) = x$ and $\gamma(1) = y$ while $\gamma(t) \in S$ for all $t \in [0, 1]$. A function $f : S \rightarrow \mathbb{R}$ is geodesically convex if S is geodesically convex and for all geodesic segments $\gamma : [0, 1] \rightarrow M$ such that $\gamma(0) = x$, $\gamma(1) = y$ and $\gamma(t) \in S$ for all $t \in [0, 1]$, the map $f \circ \gamma : [0, 1] \rightarrow \mathbb{R}$ is convex.

We recall the following lemma whose proof is standard:

LEMMA 1. *Let \mathcal{M} be a Riemannian manifold, and let $f : \mathcal{M} \rightarrow \mathbb{R}$ be a real function. Then f is geodesically convex iff Hessian of f , $\text{Hess}f$, is semi-positive definite. Also for C^2 maps $\psi, \phi : \mathcal{M} \rightarrow \mathbb{R}$, and a C^1 vector field v , we have*

$$\langle \nabla \langle \nabla \psi, \nabla \phi \rangle, v \rangle = \text{Hess}\phi(\nabla \psi, v) + \text{Hess}\psi(\nabla \phi, v). \quad (1)$$

The following basic observation for convex maps defined on a Riemannian manifold \mathcal{M} motivates for our definition for convexity of maps defined on discrete spaces.

PROPOSITION 1. *Let \mathcal{M} be a geodesically complete Riemannian manifold, and let*

$$f : \mathcal{M} \rightarrow \mathbb{R},$$

be a differentiable real function with compact support. The functional \mathcal{L}^f defined by

$$\mathcal{L}^f : \mathcal{P}_2(\mathcal{M}) \rightarrow \mathbb{R}, \quad \mathcal{L}^f(\mu) = \int_{\mathcal{M}} f d\mu.$$

is geodesically convex on $\mathcal{P}_2(\mathcal{M})$ iff, $\text{Hess}(f)$ is semi-positive definite every where on \mathcal{M} .

Proof. Here we provide an informal proof without treating regularity which makes use of Hessian. We avoid regularity discussion and we refer to [7] for a rigorous proof using MacCann-Brenier theorem. Assume that \mathcal{L}^f is geodesically convex on $\mathcal{P}_2(\mathcal{M})$. For a geodesic in $\mathcal{P}_2(\mathcal{M})$ given by the continuity equation

$$\begin{cases} \frac{\partial}{\partial t} \mu_t + \nabla \cdot (\mu_t \nabla \psi_t) = 0, \\ \frac{\partial}{\partial t} \psi_t + \frac{|\nabla \psi_t|^2}{2} = 0, \end{cases}$$

we have

$$\frac{d}{dt} \mathcal{L}^f(\mu_t) = \frac{d}{dt} \int f(x) d\mu_t = \int \langle \nabla f, \nabla \psi_t \rangle d\mu_t.$$

By applying lemma 1, we get

$$\begin{aligned}
\frac{d^2}{dt^2} \mathcal{L}^f(\mu_t) &= \frac{d}{dt} \left[\int \langle \nabla f(x), \nabla \psi_t \rangle d\mu_t \right] \\
&= \int \left\langle \nabla f, -\nabla \frac{|\nabla \psi_t|^2}{2} \right\rangle d\mu_t + \int \langle \nabla \langle \nabla f, \nabla \psi_t \rangle, \nabla \psi_t \rangle d\mu_t \\
&= \frac{-1}{2} \int \langle \nabla \langle \nabla \psi_t, \nabla \psi_t \rangle, \nabla f \rangle d\mu_t + \int \langle \nabla \langle \nabla f, \nabla \psi_t \rangle, \nabla \psi_t \rangle d\mu_t \\
&\stackrel{(1)}{=} \frac{-1}{2} \int 2 \text{Hess} \psi(\nabla \psi_t, \nabla f) d\mu_t + \int \text{Hess} f(\nabla \psi_t, \nabla \psi_t) d\mu_t \\
&\quad + \int \text{Hess} \psi(\nabla f, \nabla \psi_t) d\mu_t = \int \text{Hess} f(\nabla \psi_t, \nabla \psi_t) d\mu_t. \tag{2}
\end{aligned}$$

Hence, if f is convex, one can deduce that

$$\frac{d^2}{dt^2} \mathcal{L}^f(\mu_t) = \int \text{Hess} f(\nabla \psi_t, \nabla \psi_t) d\mu_t \geq 0.$$

Conversely, assume that \mathcal{L}^f is geodesically convex. Consider a geodesic in $\mathcal{P}_2(\mathcal{M})$ with the initial condition $(\mu_0, \nabla \psi_0)$, where $\mu_0 = \delta_{x_0}$, $\nabla \psi_0 = v$, $x_0 \in \mathcal{M}$, and $v \in T_{x_0} \mathcal{M}$. Then from (2), it can be seen that

$$0 \leq \int \text{Hess} f(\nabla \psi_0, \nabla \psi_0) d\mu_0 = \int \text{Hess} f(v, v) d\delta_{x_0}(y) = \text{Hess} f(x_0)(v, v).$$

This means that $\text{Hess} f$ is semi-positive definite. □

Inversely, this theorem reduces the study of the convexity of \mathcal{L}^f on $\mathcal{P}_2(\mathcal{M})$ to that of the Hessian of the map f at all the points of \mathcal{M} . In order to carry out the same procedure in discrete case we employ the theory developed by Maas and Erbar for discrete finite Markov spaces [8].

They consider an irreducible Markov kernel $Q : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$, on a finite set \mathcal{X} . Satisfying the condition that for every $x, y \in \mathcal{X}$, there exists a sequence $\{x_i\}_{i=0}^n \in \mathcal{X}$, such that

$$x_0 = x, \quad x_n = y, \quad \text{and} \quad Q(x_{i-1}, x_i) > 0 \quad \forall i \in \{1, 2, \dots, n\}.$$

The unique stationary probability measure π on \mathcal{X} associated with the Markov kernel Q , is assumed to be reversible, i.e. the detailed balance equation holds:

$$Q(x, y)\pi(x) = Q(y, x)\pi(y) \quad \forall x, y \in \mathcal{X}.$$

Maas and Erbar define the space of all probability densities on \mathcal{X} with respect to π as follows:

$$\mathcal{P}_\pi(\mathcal{X}) = \left\{ \rho : \mathcal{X} \rightarrow \mathbb{R}_+ \mid \sum_{x \in \mathcal{X}} \pi(x) \rho(x) = 1 \right\}. \tag{3}$$

For $\rho_0, \rho_1 \in \mathcal{P}_\pi(\mathcal{X})$, Wasserstein-Mass metric \mathcal{W} on $\mathcal{P}_\pi(\mathcal{X})$ is defined by [9]

$$\mathcal{W}(\rho_0, \rho_1)^2 = \inf_{\rho, \psi} \left\{ \frac{1}{2} \int_0^1 \sum_{x, y \in \mathcal{X}} (\psi_t(x) - \psi_t(y))^2 \Lambda(\rho_t(x), \rho_t(y)) Q(x, y) \pi(x) dt \right\}, \tag{4}$$

where the infimum is taken over all piece-wise C^1 curves $\rho : [0, 1] \rightarrow \mathcal{P}_\pi(\mathcal{X})$, and the functions $\psi : [0, 1] \times \mathcal{X} \rightarrow \mathbb{R}$ satisfies a certain continuity equation. The map Λ is taken to be an arithmetic average for our case. The space $(\mathcal{P}_\pi(\mathcal{X}), \mathcal{W})$ has the structure of a complete Riemannian manifold.

PROPOSITION 2 [8]. For any $\bar{\rho} \in \mathcal{P}_\pi(\chi)$ and $\bar{\psi} \in \mathbb{R}^\chi$, on a sufficiently small interval around 0, there exists a unique constant speed geodesic such that $\rho_0 = \bar{\rho}$ and initial tangent vector $\nabla \psi_0 = \nabla \bar{\psi}$ satisfying

$$\begin{cases} \partial_t \rho_t(x) = - \sum_{y \in \chi} (\psi_t(y) - \psi_t(x)) \Lambda(\rho_t(x), \rho_t(y)) Q(x, y), \\ \partial_t \psi_t(x) = - \frac{1}{2} \sum_{y \in \chi} (\psi_t(x) - \psi_t(y))^2 \partial_1 \Lambda(\rho_t(x), \rho_t(y)) Q(x, y). \end{cases} \quad (5)$$

Definition of convex functions on Markov spaces. Let (χ, Q, π) be a Markov triple on a finite space χ as described above. Let $f : \chi \rightarrow \mathbb{R}$ be a real valued function on χ . Inspired by the observation given in Proposition 1, we want to introduce a notion of convexity for f . To this end, consider the functional L_π^f

$$\begin{aligned} L_\pi^f : \mathcal{P}_\pi(\chi) &\longrightarrow \mathbb{R} \\ L_\pi^f(\rho) &= \sum_{w \in \mathcal{X}} f(w) \rho(w). \end{aligned} \quad (6)$$

In the discrete case, we are interested in geodesic convexity of L_π^f on a totally geodesic submanifold of $\mathcal{P}_\pi(\chi)$ passing through all the vertices of the simplex $\mathcal{P}_\pi(\chi)$ (see Definition 1). In the case where χ has only two points, $\mathcal{P}_\pi(\chi)$ will become a 1-simplex in \mathbb{R}^2 . We start by identifying geodesically convex functionals L_π^f on this 1-simplex.

Let us denote by $\mathcal{P}(\chi)$, the space of probability measures on χ , which consists of all maps $\hat{\rho} : \chi \rightarrow [0, 1]$ satisfying $\sum_{x \in \chi} \hat{\rho}(x) = 1$. Let \mathcal{H} be the map defined by

$$\mathcal{H} : \mathcal{P}_\pi(\chi) \longrightarrow \mathcal{P}(\chi), \quad (\mathcal{H}(\rho))(x) = \pi(x) \rho(x).$$

Since \mathcal{H} is a bijection map, we can use it to equip $\mathcal{P}(\chi)$ with a Riemannian metric \mathcal{G} with respect to which \mathcal{H} is an isometry. The Riemannian metric on $\mathcal{P}_\pi(\chi)$ is defined in [8].

Given a real function $f : \chi \rightarrow \mathbb{R}$, let $L^f : \mathcal{P}(\chi) \rightarrow \mathbb{R}$ be defined by

$$L^f(\hat{\rho}) = \sum_{x \in \chi} \hat{\rho}(x) f(x).$$

Our aim is to study the geodesic convexity of L^f with respect to \mathcal{G} .

Based on the observation of Proposition 1, we introduce a notion of convexity for real functions $f : \chi \rightarrow \mathbb{R}$ on the discrete Markov space (χ, Q, π) . We denote by \mathcal{G}_π the Erbar-Maas metric on $\mathcal{P}_\pi(\chi)$, and we set $\mathcal{G} := ((\mathcal{H})^{-1})^* \mathcal{G}_\pi$, which is a Riemannian metric on $\mathcal{P}(\chi)$.

Definition 1. For a real $\lambda \in \mathbb{R}$, a map $f : \chi \rightarrow \mathbb{R}$ is called λ -convex, if there exists a totally geodesic submanifold $\mathcal{A} \subset \mathcal{P}(\chi)$ containing all the vertices of the simplex $\mathcal{P}(\chi)$ such that $L^f : \mathcal{A} \rightarrow \mathbb{R}$ is λ -convex with respect to \mathcal{G} . This is equivalent to say that $L_\pi^{\pi f} : \mathcal{A}_\pi \rightarrow \mathbb{R}$ to be λ -convex with respect to \mathcal{G}_π , where $\mathcal{A}_\pi \subset \mathcal{P}_\pi(\chi)$ is a totally geodesic submanifold containing all the vertices of $\mathcal{P}_\pi(\chi)$.

2.1. Tensorisation and hypercube structure

Let (χ_i, Q_i, π_i) , for $1 \leq i \leq n$, be irreducible and reversible Markov chains. Let $Q_{(i)}$ denote the lift of Q_i to the product space $\chi = \chi_1 \times \cdots \times \chi_n$, defined for $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_n)$ by

$$Q_{(i)}(x, y) = \begin{cases} Q_i(x_i, y_i), & \text{If } x_j = y_j \text{ for all } j \neq i, \\ 0, & \text{Otherwise.} \end{cases}$$

For $\alpha = (\alpha_1, \dots, \alpha_n)$ satisfying $\sum_{i=1}^n \alpha_i = 1$, $\alpha_i \geq 0$, $1 \leq i \leq n$, consider a weighted product Markov chain

$(\mathcal{X}, Q_\alpha, \pi)$, such that

$$Q_\alpha := \sum_{i=1}^n \alpha_i Q_{(i)},$$

and where the corresponding reversible probability measure is given by $\pi = \pi_1 \otimes \cdots \otimes \pi_n$.

Let (\mathcal{X}, Q, π) denote the above tensorized Markov space where we are setting $Q := Q_\alpha$ for $\alpha = (\frac{1}{n}, \dots, \frac{1}{n})$, therefore

$$Q(x, y) = \sum_{i=1}^n \frac{1}{n} Q_{(i)}(x, y). \quad (7)$$

Let $\rho_t^i \in \mathcal{P}_{\pi_i}(\mathcal{X}_i)$ be a geodesic and for any $x = (x_1, \dots, x_n) \in \mathcal{X}$, we define

$$\rho_t(x) = \left(\rho_{\frac{t}{n}}^1 \times \cdots \times \rho_{\frac{t}{n}}^n \right) (x) := \prod_{i=1}^n \rho_{\frac{t}{n}}^i(x_i). \quad (8)$$

It is not difficult to see that $\rho_t \in \mathcal{P}_\pi(\mathcal{X})$.

Also for ψ_t^i satisfying the geodesic equation (5), we set

$$\psi_t(x) := \sum_{i=1}^n \psi_{\frac{t}{n}}^i(x_i). \quad (9)$$

Our claim is that $t \rightarrow \rho_t$ is a geodesic in $\mathcal{P}_\pi(\mathcal{X})$:

THEOREM 1. *The pair (ρ_t, ψ_t) satisfies the system of equations (5), which means that $t \rightarrow \rho_t$ defines a geodesic in $\mathcal{P}_\pi(\mathcal{X})$.*

Proof. From (8) and (5) we have

$$\begin{aligned} \frac{d}{dt} \rho_t(x) &= \frac{d}{dt} \left(\rho_{\frac{t}{n}}^1 \times \cdots \times \rho_{\frac{t}{n}}^n \right) (x_1, \dots, x_n) = \frac{d}{dt} \prod_{i=1}^n \rho_{\frac{t}{n}}^i(x_i) = \sum_{j=1}^n \frac{d}{dt} \rho_{\frac{t}{n}}^j(x_j) \prod_{i \neq j} \rho_{\frac{t}{n}}^i(x_i) \\ &= -\frac{1}{n} \sum_{j=1}^n \prod_{i \neq j} \rho_{\frac{t}{n}}^i(x_i) \sum_{y_j \in \mathcal{X}_j} (\psi_{\frac{t}{n}}^j(y_j) - \psi_{\frac{t}{n}}^j(x_j)) \Lambda(\rho_{\frac{t}{n}}^j(y_j), \rho_{\frac{t}{n}}^j(x_j)) Q_j(x_j, y_j) \\ &= -\frac{1}{n} \sum_{j=1}^n \sum_{y_j \in \mathcal{X}_j} (\psi_{\frac{t}{n}}^j(y_j) - \psi_{\frac{t}{n}}^j(x_j)) \Lambda(\rho_{\frac{t}{n}}^j(y_j) \prod_{i \neq j} \rho_{\frac{t}{n}}^i(x_i), \rho_{\frac{t}{n}}^j(x_j) \prod_{i \neq j} \rho_{\frac{t}{n}}^i(x_i)) Q_j(x_j, y_j) \\ &= -\sum_{y \in \mathcal{X}, y \sim x} (\psi_t(y) - \psi_t(x)) \Lambda(\rho_t(y), \rho_t(x)) Q(x, y). \end{aligned} \quad (10)$$

Here, in the third line, we used property $\Lambda(lt, ls) = l\Lambda(t, s)$, for $l > 0$ and $s, t \geq 0$, and, in the last line, we used the definition of ψ_t in (9) and Q in (7). Also, $y \sim x$ means y is a neighborhood of x with respect to the hypercube structure over \mathcal{X} .

From (9), it can also be seen that

$$\begin{aligned} \frac{d}{dt} \psi_t &= -\frac{1}{2n} \sum_{i=1}^n \sum_{y_i \in \mathcal{X}_i} \left(\psi_{\frac{t}{n}}^i(y_i) - \psi_{\frac{t}{n}}^i(x_i) \right)^2 \partial_1 \Lambda \left(\rho_{\frac{t}{n}}^i(x_i), \rho_{\frac{t}{n}}^i(y_i) \right) Q_i(x_i, y_i) \\ &= -\frac{1}{2} \sum_{y \sim x} (\psi_t(y) - \psi_t(x))^2 \partial_1 \Lambda \left(\prod_{i=1}^n \rho_{\frac{t}{n}}^i(x_i), \prod_{i=1}^n \rho_{\frac{t}{n}}^i(y_i) \right) Q(x, y). \end{aligned} \quad (11)$$

(10) and (11) are nothing but the geodesic equations for the pair (ρ_t, ψ_t) on $\mathcal{P}_\pi(\mathcal{X})$. Note that here again, we used the derivative of the relation $\Lambda(lt, ls) = l\Lambda(t, s)$, for $l > 0$ and $s, t \geq 0$, and the definition of Q .

Consider the spaces $\mathcal{A} \subset \mathcal{P}(\chi)$ and $\mathcal{A}_\pi \subset \mathcal{P}_\pi(\chi)$ respectively defined by

$$\mathcal{A} = \left\{ \prod_{i=1}^n \rho_i \mid \rho_i \in \mathcal{P}(\chi_i) \right\}, \quad \mathcal{A}_\pi = \left\{ \prod_{i=1}^n \hat{\rho}_i \mid \hat{\rho}_i \in \mathcal{P}_\pi(\chi_i) \right\}.$$

COROLLARY 1. *The spaces \mathcal{A} and \mathcal{A}_π are totally geodesic subspaces of $\mathcal{P}(\chi)$ and $\mathcal{P}_\pi(\chi)$, respectively.*

3. HESSIAN OF A REAL MAP DEFINED ON HYPERCUBE

Now consider n two-point Markov chains (χ_k, Q_k, π_k) , where as before $\chi_k = \{\pm 1\}$, for $k = 1, \dots, n$. Let (χ, Q, π) denote the tensorized product space. If for $k = 1, \dots, n$ we set $U := Q_k(+1, -1)$, and $V := Q_k(-1, +1)$, then by stationarity condition (3) $\pi_k(+1) = \frac{V}{U+V}$, and $\pi_k(-1) = \frac{U}{U+V}$. We assume that $U > V$, also Λ is assumed to be arithmetic mean then from the continuity equation (10) one can deduce that

$$\frac{d}{dt} L^f(\rho_t) = \sum_x \Gamma(f, \psi_t)(x) \rho_t(x) \pi(x)$$

where $\Gamma(f, g)(x) := \frac{1}{2} \sum_y Q(x, y) (f(x) - f(y)) (g(x) - g(y))$. A second derivation along a geodesic ρ_t leads to

$$\frac{d^2}{dt^2} L^f(\rho_t) = \sum_x \left(\Gamma \left(\Gamma(f, \psi_t), \psi_t \right) - \frac{1}{2} \Gamma(f, \Gamma(\psi_t)) \right) (x) \rho_t(x) \pi(x)$$

Hence if we define the Hessian of f as follows

$$\text{Hess}(f)[\nabla \phi, \nabla \psi] = \Gamma(\Gamma(f, \phi), \psi) - \frac{1}{2} \Gamma(f, \Gamma(\phi, \psi))$$

then by taking $\rho_t(x) = \delta_{x_0}$ it can be seen that L^f is geodesically convex iff $\text{Hess}(f)$ is positive definite at all $x \in \chi$. We have

$$\begin{aligned} \Gamma(\Gamma(f, \phi), \psi)(x) &= \frac{1}{2} \sum_y Q(x, y) (\Gamma(f, \phi)(x) - \Gamma(f, \phi)(y)) (\psi(x) - \psi(y)) \\ &= \frac{1}{2} \sum_y Q(x, y) \left(\frac{1}{2} \sum_z Q(x, z) (f(x) - f(z)) (\phi(x) - \phi(z)) \right. \\ &\quad \left. - \frac{1}{2} \sum_z Q(y, z) (f(y) - f(z)) (\phi(y) - \phi(z)) \right) (\psi(x) - \psi(y)) \\ &= \frac{1}{4} \sum_{y, z} Q(x, y) Q(x, z) (f(x) - f(z)) (\phi(x) - \phi(z)) (\psi(x) - \psi(y)) \\ &\quad - \frac{1}{4} \sum_{y, z} Q(x, y) Q(y, z) (f(y) - f(z)) (\phi(y) - \phi(z)) (\psi(x) - \psi(y)) \end{aligned}$$

and

$$\begin{aligned} \Gamma(\Gamma(\phi, \psi), f)(x) &= \frac{1}{4} \sum_{y, z} Q(x, y) Q(x, z) (\phi(x) - \phi(z)) (\psi(x) - \psi(y)) (f(x) - f(y)) \\ &\quad - \frac{1}{4} \sum_{y, z} Q(x, y) Q(y, z') (\phi(y) - \phi(z')) (\psi(y) - \psi(z')) (f(x) - f(y)) \end{aligned}$$

Now at a fixed $x \in \chi$ and for $k = 1, \dots, n$ we define $\phi_k : \chi \rightarrow \mathbb{R}$ by $\phi_k(y) = \begin{cases} 0 & \text{if } y_k = x_k \\ 1 & \text{Otherwise} \end{cases}$. Then if r_k

for $k = 1, \dots, n$ denotes the operator $r_k(\theta_1, \dots, \theta_k, \dots, \theta_n) = (\theta_1, \dots, -\theta_k, \dots, \theta_n)$, then since $(\phi(r_k(x)) - \phi(r_l \circ r_k(x))) (\psi(r_k(x)) - \psi(r_l \circ r_k(x))) = 0$ we obtain

$$\begin{aligned} & [\text{Hess}(f)[\nabla\phi_k, \nabla\phi_l](x) + \text{Hess}(f)[\nabla\phi_l, \nabla\phi_k](x) \\ &= \left[\frac{1}{8} \mathcal{Q}(x, r_k(x)) \mathcal{Q}(x, r_l(x)) (2f(x) - f \circ r_k(x) - f \circ r_l(x)) \right. \\ & \quad - \frac{1}{4} \mathcal{Q}(x, r_k(x)) \mathcal{Q}(r_k(x), r_l \circ r_k(x)) (f(r_k(x)) - f(r_l \circ r_k(x))) \\ & \quad \left. - \frac{1}{4} \mathcal{Q}(x, r_l(x)) \mathcal{Q}(r_l(x), r_k \circ r_l(x)) (f(r_l(x)) - f(r_k \circ r_l(x))) \right] \\ &= \frac{UV}{n^2} \left[\frac{1}{4} f(x) - \frac{3}{8} f(r_k(x)) - \frac{3}{8} f(r_l(x)) + \frac{1}{2} f(r_k \circ r_l(x)) \right] \end{aligned}$$

and

$$\begin{aligned} \text{Hess}(f)[\nabla\phi_k, \nabla\phi_k](x) &= \frac{1}{4} \mathcal{Q}(x, r_k(x)) \mathcal{Q}(x, r_k(x)) (f(x) - f(r_k(x))) \\ & \quad - \frac{1}{4} \mathcal{Q}(r_k(x), x) \mathcal{Q}(r_k(x), x) (f(x) - f(r_k(x))) \\ &= \frac{1}{\pi_k(x_k)} \frac{UV(U-V)}{4n^2(U+V)} \left[f(x_1, \dots, x_{k-1}, +1, x_{k+1}, \dots, x_n) \right. \\ & \quad \left. - f(x_1, \dots, x_{k-1}, -1, x_{k+1}, \dots, x_n) \right] \end{aligned}$$

We set $H_{kl} = \frac{1}{2} [\text{Hess}(f)[\nabla\phi_k, \nabla\phi_l](x) + \text{Hess}(f)[\nabla\phi_l, \nabla\phi_k](x)$ and we define the matrix $H := [H_{k,l}]_{1 \leq k, l \leq n}$.

Definition 2. We call H the Hessian of the map f at the point $(\theta_1, \dots, \theta_n)$.

We have proved that

THEOREM 2. *The map $f : \chi \rightarrow \mathbb{R}$ is strictly convex If H is positive definite at all $(\theta_1, \dots, \theta_n)$.*

THEOREM 3. *Given a strictly convex function $f : \chi \rightarrow \mathbb{R}$ on a Markov chain (χ, \mathcal{Q}, π) in the sense of Definition 1, then f admits a unique minimum on χ .*

Proof. L^f is linear on the simplex $\mathcal{P}(\chi)$. Hence the minimum of L^f occurs at a unique vertex of the simplex $\mathcal{P}(\chi)$. On the other hand, by definition 1 we know that $L^f|_{\mathcal{A}}$ is geodesically strictly convex. So the minimum of L^f can only occur on one of the vertices of \mathcal{A} , and moreover it is unique. The vertices of $\mathcal{P}(\chi)$ are of the form δ_x for some $x \in \chi$, where δ_x denotes the Dirac probability distribution on χ concentrated at x . Moreover, by definition of L^f we know that $L^f(\delta_x) = f(x)$. This proves that the minimum of f can only occur at a single point in χ . □

3.1. Examples of convex maps on hypercube and the Ising model

Consider a linear map $f : \chi \rightarrow \mathbb{R}$ given by the following relation

$$f(\theta_1, \dots, \theta_n) = \sum_{k=1}^n a_k \theta_k$$

Then the Hessian of f at all the points $(\theta_1, \dots, \theta_n)$ is diagonal and it is positive definite exactly when the diagonal coefficients are positive. It is not difficult to see that this occurs when $a_k > 0$ for all k .

The quadratic map $f : \chi \rightarrow \mathbb{R}$ defined by

$$f(\theta_1, \dots, \theta_n) = \sum_{k < l} J_{kl} \theta_k \theta_l,$$

is never convex. To see this, we observe that for $k \neq l$ $H_{kl} = 4J_{kl}$ and

$$H_{kk} = f(\theta_1, \dots, \theta_{k-1}, 1, \theta_{k+1}, \dots, \theta_n) - f(\theta_1, \dots, \theta_{k-1}, -1, \theta_{k+1}, \dots, \theta_n) = 2 \sum_{j \neq k} J_{kj} \theta_j.$$

There always exists choices for $\theta_1, \dots, \theta_n$, for which the above quantity H_{kk} is negative.

As is explained in [10], mean field equations in statistical physics can be obtained by variational method applied to the Kullback-Leibler divergence, where $\mathcal{A} \subset \mathcal{P}(\chi)$ consists of the space of factorized probability distributions. More precisely in an Ising model with the energy function $\mathcal{H}(\theta_1, \dots, \theta_n) = \sum_{k < l} J_{kl} \theta_k \theta_l + \sum_{k=1}^n a_k \theta_k$, where J_{kl} s are the coupling constants representing the strength of interaction between the atoms k and l and a_k 's the external field of the system, the so called mean field equation can be derived through a variational problem of the form $\min_{q \in \mathcal{A}} D_{KL}(q||p)$. Now due to the fact that hyper cubes' Ricci curvature is non-negative according to [9] one can derive the following theorem,

THEOREM 3. *If the energy function $\mathcal{H}(\theta_1, \dots, \theta_n) = \sum_{k < l} J_{kl} \theta_k \theta_l + \sum_{k=1}^n a_k \theta_k$ is convex, then the $D_{KL}(\cdot||p) : \mathcal{A} \rightarrow \mathbb{R}$ is convex, and therefore the optimization problem $\min_{q \in \mathcal{A}} D_{KL}(q||p)$ admits a unique solution. Here p denote the probability distribution $p(\theta_1, \dots, \theta_n) = \frac{e^{\mathcal{H}(\theta_1, \dots, \theta_n)}}{\sum_{\theta} e^{\mathcal{H}(\theta_1, \dots, \theta_n)}}$.*

ACKNOWLEDGEMENTS

The author would like to thank the anonymous referee for valuable comments and suggestions on this paper.

REFERENCES

- [1] Bahraini A. Mean field variational Bayesian inference for Gaussian mixture model rigorous uncertainty quantification. Submitted. 2024.
- [2] Jordan R, Kinderlehrer D, Otto F. The variational formulation of the Fokker-Planck equation. SIAM J. Math. Anal. 1998; 29(1): 1–17. DOI: 10.1137/S0036141096303359.
- [3] Otto F. The geometry of dissipative evolution equations: the porous medium equation. Comm. Partial Differential Equations. 2001; 26(1-2): 101–174. DOI: 10.1081/PDE-100002243.
- [4] Ambrosio L, Gigli N, Savare G. Gradient flows in metric spaces and in the Wasserstein space of probability measures. Lectures in Mathematics. ETH Zürich: Birkhäuser Basel; 2005. DOI: 10.1007/b137080.
- [5] Blei DM, Kucukelbir A, McAuliffe JD, Variational inference: a review for statisticians. J. Amer. Statist. Assoc. 2017; 112(518): 859–877.
- [6] Arjovsky M, Chintala S, Bottou L. Wasserstein generative adversarial networks. In: Precup D, Teh YW. Proceedings of the 34th International Conference on Machine Learning. Proc. of Machine Learning Research (PMLR), vol. 70. 2017, pp. 214–223.
- [7] Bahraini A, Sadeghi S. Optimal transport and variational Bayesian inference. International Journal of Approximate Reasoning. 2023; 162: 109022.
- [8] Erbar M, Maas J. Ricci curvature of finite Markov chains via convexity of the entropy. Arch. Rational. Mech. Anal. 2012; 206: 997–1038. DOI: 10.1007/s00205-012-0554-z.
- [9] Maas J. Gradient flows of the entropy for finite Markov chains. Journal of Functional Analysis. 2011; 261(8): 2250–2292. DOI: 10.1016/j.jfa.2011.06.009.
- [10] Opper M, Winther O. From naive mean field theory to the TAP equations. In: Advanced mean field methods: theory and practice. MIT Press; 2001, ch. 2, pp. 7–20.
- [11] Lott J, Villani C. Ricci curvature for metric-measure spaces via optimal transport. Ann. of Math. 2009; 169(3): 903–991.
- [12] Otto F. The geometry of dissipative evolution equations: the porous medium equation. Comm. Partial Differential Equations. 2001; 26(1-2): 101–174. DOI: 10.1081/PDE-100002243.

Received September 22, 2024