

Wasserstein Space

Lecturer: Jiaming Liang

April 4, 2024

1 Wasserstein space

Recall from Lecture 8 that the Fokker-Planck equation is

$$\frac{\partial \rho_t}{\partial t} = \nabla \cdot \left(\rho_t \nabla \log \frac{\rho_t}{\nu} \right) = \nabla \cdot (\rho_t \nabla f) + \Delta \rho_t. \quad (1)$$

We also claimed that a variational perspective on (1) is that it is the gradient flow for minimizing the relative entropy $\text{KL}(\cdot \| \nu)$ in the space of probability distributions with the Wasserstein-2 metric. The objective of this lecture is to develop necessary calculus to substantiate this claim, understand the Wasserstein space, and analyze convergence along the Wasserstein gradient flow.

Background on Riemannian geometry A manifold \mathcal{M} is a space that is locally homeomorphic to a Euclidean space. The tangent space $T_p \mathcal{M}$ at $p \in \mathcal{M}$ is an associated vector space containing all possible velocities (tangent bundle) of curves $\gamma(t)$ passing through p . A Riemannian metric is a choice of inner products $p \mapsto \langle \cdot, \cdot \rangle_p$ on the tangent spaces, depending smoothly on p .

The Riemannian metric induces a distance function (in the sense of metric spaces) via

$$d(p, q) := \inf \left\{ \int_0^1 \|\dot{\gamma}(t)\|_{\gamma(t)} dt \mid \gamma : [0, 1] \rightarrow \mathcal{M}, \gamma(0) = p, \gamma(1) = q \right\}, \quad (2)$$

where $\dot{\gamma}(t) \in T_{\gamma(t)} \mathcal{M}$ denotes the tangent vector and $\|\cdot\|_{\gamma(t)}$ is induced by $\langle \cdot, \cdot \rangle_{\gamma(t)}$. If the infimum is achieved by a curve γ , then γ is referred to as a geodesic (a shortest path); if $t \mapsto \|\dot{\gamma}(t)\|_{\gamma(t)}$ is constant, then it is called a constant-speed geodesic. We say (\mathcal{M}, d) is a Riemannian manifold.

Given a functional $\mathcal{F} : \mathcal{M} \rightarrow \mathbb{R}$, the gradient of \mathcal{F} at p is defined to be the unique element $\nabla \mathcal{F}(p) \in T_p \mathcal{M}$ such that for all curves $(\rho(t))_{t \in \mathbb{R}}$ passing through p at time 0 with velocity $v \in T_p \mathcal{M}$ (i.e., $\rho(0) = p$ and $\rho'(0) = v$), it holds that $\frac{d}{dt} \mathcal{F}(\rho(t))|_{t=0} = \langle \nabla \mathcal{F}(p), v \rangle_p$.

In the rest of the lecture, we consider a specific Riemannian manifold, namely the Wasserstein space $(\mathcal{P}_2(\mathbb{R}^d), \mathcal{W}_2)$, and develop the differential calculus in the Wasserstein space.

To develop the analogy of \mathcal{W}_2 and metric, we would like to define a metric structure $\langle \cdot, \cdot \rangle_\rho$ on each tangent space $T_\rho \mathcal{P}_2$, depending smoothly on ρ . This metric should define a norm $\|\cdot\|_\rho$ on each $T_\rho \mathcal{P}_2$ such that

$$\mathcal{W}_2(\rho_0, \rho_1)^2 = \inf \left\{ \int_0^1 \left\| \frac{\partial \rho}{\partial t} \right\|_{\rho(t)}^2 dt \mid \rho : [0, 1] \rightarrow \mathcal{P}_2(\mathbb{R}^d), \rho(0) = \rho_0, \rho(1) = \rho_1 \right\},$$

where the infimum is taken over all paths connecting ρ_0 and ρ_1 .

A tangent vector $v \in T_\rho \mathcal{P}_2$ is $\frac{d\rho}{dt}|_{t=0}$ where $\rho(t)$ satisfies $\rho(0) = \rho$. From a fluid dynamics perspective, we want to see the path $\rho(t)$ as the time-evolving density of a set of particles moving continuously from ρ_0 to ρ_1 with a velocity field $c_t(x)$, that is the velocity of a particle is uniquely determined by its position. Now, we have the dynamics of the particles

$$\dot{x}_t = c_t(x_t). \quad (3)$$

We want to study the evolution of their law $\rho_t(x) = \rho(t, x)$.

Lemma 1. *The continuity equation of (3) is*

$$\frac{\partial \rho_t(x)}{\partial t} + \nabla \cdot (\rho_t(x) \nabla \lambda(x)) = 0$$

for some smooth function λ . The length of the tangent vector is given by

$$\left\| \frac{\partial \rho}{\partial t} \right\|_{\rho(t)}^2 = \mathbb{E}_{X \sim \rho_t} [\|\nabla \lambda(X)\|^2].$$

Proof. It follows from Theorem 1 of Lecture 8 that the continuity equation is

$$\frac{\partial \rho_t}{\partial t} + \nabla \cdot (\rho_t c_t) = 0.$$

The velocity field c_t is not unique: e.g., let w be a vector field with zero divergence (i.e., $\nabla \cdot w = 0$), then $c'_t = c_t + \varepsilon w / \rho_t$ for $\varepsilon \neq 0$ is also admissible since

$$\frac{\partial \rho_t(x)}{\partial t} + \nabla \cdot \left[\rho_t \left(c_t + \varepsilon \frac{w}{\rho_t} \right) \right] = 0. \quad (4)$$

Among all possible vector fields, we want to select the one with lowest kinetic energy, if possible. That is

$$\begin{aligned} & \inf_{c_t} \frac{1}{2} \int \rho_t(x) \|c_t(x)\|^2 dx \\ & \text{s.t. } -\nabla \cdot (\rho_t(x) c_t(x)) = \frac{\partial \rho_t(x)}{\partial t} \end{aligned}$$

Approach 1: Consider the Lagrange function with multiplier $\lambda(x)$

$$\begin{aligned} & \inf_{c_t} \frac{1}{2} \int \rho_t(x) \|c_t(x)\|^2 dx + \int \lambda(x) \nabla \cdot (\rho_t(x) c_t(x)) dx \\ & = \inf_{c_t} \frac{1}{2} \int \rho_t(x) \|c_t(x)\|^2 dx - \int \langle \nabla \lambda(x), c_t(x) \rangle \rho_t(x) dx. \end{aligned}$$

The problem boils down to a pointwise minimization and gives

$$c_t(x) = \nabla\lambda(x). \quad (5)$$

Approach 2: As in (4), $c_t + \varepsilon w/\rho_t$ is a feasible solution. Assume that c_t is a minimizing vector field, then we have

$$\int \rho_t \|c_t\|^2 \leq \int \rho_t \left\| c_t + \varepsilon \frac{w}{\rho_t} \right\|^2$$

Using an limiting argument (i.e., $\varepsilon \rightarrow 0$), we can show

$$0 = \int \langle c_t, w \rangle,$$

that is c_t is orthoganal to the set of divergence-free vector fields. This means c_t should be a gradient $\nabla\lambda$, so that

$$\int \langle c_t, w \rangle = \int \langle \nabla\lambda, w \rangle = \int \lambda \nabla \cdot w = 0.$$

In both approaches, we prove that (5) holds for some smooth function λ . Moreover, the norm $\|\cdot\|_\rho$ is given by the minimum kinetic energy

$$\left\| \frac{\partial \rho_t}{\partial t} \right\|_{\rho(t)}^2 = \int \rho_t(x) \|c_t(x)\|^2 dx = \mathbb{E}_{X \sim \rho_t} [\|\nabla\lambda(X)\|^2].$$

□

A sufficient condition for uniqueness is that ρ_t satisfies the Poincaré inequality. Standard results on elliptic PDE assert the existence and uniqueness.

To sum up, one can write

$$\left\{ \begin{array}{l} \mathcal{W}_2^2(\rho_0, \rho_1) = \inf \left\{ \int_0^1 \left\| \frac{\partial \rho}{\partial t} \right\|_{\rho(t)}^2 dt \mid \rho : \rho(0) = \rho_0, \quad \rho(1) = \rho_1 \right\} \\ \left\| \frac{\partial \rho}{\partial t} \right\|_\rho^2 = \int \rho \|\nabla\lambda\|^2, \quad -\nabla \cdot (\rho \nabla\lambda) = \frac{\partial \rho}{\partial t}. \end{array} \right. \quad (6)$$

This definition formally endows $\mathcal{P}_2(\mathbb{R}^d)$ a Riemannian metric structure. It comes from a differential formulation of optimal trasport.

2 Otto calculus

By polarization and the norm $\|\cdot\|_\rho$, we can define the inner product of two tangent vectors $\frac{\partial \rho}{\partial t_1} \in T_p \mathcal{P}_2$ and $\frac{\partial \rho}{\partial t_2} \in T_p \mathcal{P}_2$: first, we know there exist smooth functions λ_1 and λ_2 satisfying

$$\frac{\partial \rho}{\partial t_1} = -\nabla \cdot (\rho \nabla\lambda_1), \quad \frac{\partial \rho}{\partial t_2} = -\nabla \cdot (\rho \nabla\lambda_2).$$

Thus, we define $\langle \cdot, \cdot \rangle_p$ as follows

$$\begin{aligned} \left\langle \frac{\partial \rho}{\partial t_1}, \frac{\partial \rho}{\partial t_2} \right\rangle_\rho &= \frac{1}{4} \left(\left\| \frac{\partial \rho}{\partial t_1} + \frac{\partial \rho}{\partial t_2} \right\|_\rho^2 - \left\| \frac{\partial \rho}{\partial t_1} - \frac{\partial \rho}{\partial t_2} \right\|_\rho^2 \right) \\ &= \frac{1}{4} \left(\int \rho \|\nabla \lambda_1 + \nabla \lambda_2\|^2 dx - \int \rho \|\nabla \lambda_1 - \nabla \lambda_2\|^2 dx \right) \\ &= \int \rho \langle \nabla \lambda_1, \nabla \lambda_2 \rangle dx. \end{aligned}$$

Given a functional $\mathcal{F} : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$, the Wasserstein gradient of \mathcal{F} at p is defined to be the unique element $\mathbf{grad}_W \mathcal{F}(p) \in T_p \mathcal{P}_2$ such that for all curves $(\rho(t))_{t \in \mathbb{R}}$ passing through p at time 0 with velocity $v \in T_p \mathcal{P}_2$ (i.e., $\rho(0) = p$ and $\rho'(0) = v$), it holds that

$$\frac{d}{dt} \mathcal{F}(\rho(t))|_{t=0} = \langle \mathbf{grad}_W \mathcal{F}(p), v \rangle_p. \quad (7)$$

Definition 1. The first variation of $\mathcal{F} : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ evaluated at $\rho \in \mathcal{P}_2(\mathbb{R}^d)$ is a function $\frac{\delta \mathcal{F}}{\delta \rho} : \mathbb{R}^d \rightarrow \mathbb{R}$ that satisfies the following for all $\nu \in \mathcal{P}_2(\mathbb{R}^d)$,

$$\lim_{\epsilon \rightarrow 0} \frac{\mathcal{F}(\rho + \epsilon(\nu - \rho)) - \mathcal{F}(\rho)}{\epsilon} = \left\langle \frac{\delta \mathcal{F}}{\delta \rho}, \nu - \rho \right\rangle_{L^2(\mathbb{R}^d)} = \int_{\mathbb{R}^d} \frac{\delta \mathcal{F}}{\delta \rho}(x) (\nu(x) - \rho(x)) dx.$$

Lemma 2. The Wasserstein gradient of \mathcal{F} is

$$\mathbf{grad}_W \mathcal{F}(p) = -\nabla \cdot \left(p \nabla \frac{\delta \mathcal{F}}{\delta p} \right). \quad (8)$$

Proof. Let $p \in \mathcal{P}_2(\mathbb{R}^d)$ be given and consider a smooth curve $\rho(t)$ satisfying $\rho(0) = p$ and $\rho'(0) = v$ for some $v \in T_p \mathcal{P}_2$. Since $\mathbf{grad}_W \mathcal{F}(p) \in T_p \mathcal{P}_2$, we know there exist smooth functions f and λ solving the following continuity equations

$$\mathbf{grad}_W \mathcal{F}(p) = -\nabla \cdot (p \nabla f), \quad v = -\nabla \cdot (p \nabla \lambda). \quad (9)$$

On the one hand, by the definition of inner production $\langle \cdot, \cdot \rangle_p$, we have

$$\langle \mathbf{grad}_W \mathcal{F}(p), v \rangle_p = \int \langle \nabla f, \nabla \lambda \rangle p dx.$$

On the other hand, by definitions and integration by parts, we have

$$\begin{aligned} \frac{d}{dt} \mathcal{F}(\rho(t))|_{t=0} &= \lim_{t \rightarrow 0} \frac{\mathcal{F}(p + tv) - \mathcal{F}(p)}{t} = \left\langle \frac{\delta \mathcal{F}}{\delta p}, v \right\rangle_{L^2} \\ &= \int \frac{\delta \mathcal{F}(p)}{\delta p}(x) v(x) dx = \int \frac{\delta \mathcal{F}}{\delta p} [-\nabla \cdot (p \nabla \lambda)] dx \\ &= \int \left\langle \nabla \frac{\delta \mathcal{F}}{\delta p}, \nabla \lambda \right\rangle p dx. \end{aligned}$$

In view of the definition of Wasserstein gradient (7), comparing the above two equations and noting v (and $\nabla\lambda$) is arbitrary, we must have

$$\nabla f = \nabla \frac{\delta \mathcal{F}}{\delta \rho}.$$

Plugging the above identity into (9), we prove (8). \square

Examples of Wasserstein gradients

1) Negative entropy $\mathcal{F}(\rho) = -\mathbf{H}(\rho) = \int \rho \log \rho$,

$$\frac{\delta \mathcal{F}}{\delta \rho} = 1 + \log \rho, \quad \mathbf{grad}_W \mathcal{F}(\rho) = -\nabla \cdot \left(\rho \nabla \frac{\delta \mathcal{F}}{\delta \rho} \right) = -\nabla \cdot (\rho \nabla \log \rho) = -\Delta \rho,$$

then the Wasserstein gradient flow of entropy is the heat equation

$$\frac{\partial \rho_t}{\partial t} = -\mathbf{grad}_W \mathcal{F}(\rho_t) = \Delta \rho_t.$$

2) KL divergence (relative entropy) $\mathcal{F}(\rho) = \mathbf{KL}(\rho \parallel \nu) = \mathbf{H}_\nu(\rho) = \int \rho \log \frac{\rho}{\nu}$,

$$\frac{\delta \mathcal{F}}{\delta \rho} = 1 + \log \rho - \log \nu, \quad \mathbf{grad}_W \mathcal{F}(\rho) = -\nabla \cdot \left(\rho \nabla \frac{\delta \mathcal{F}}{\delta \rho} \right) = -\nabla \cdot (\rho \nabla \log \rho - \rho \nabla \log \nu) = -\Delta \rho - \nabla \cdot (\rho \nabla f),$$

then the Wasserstein gradient flow of KL divergence is the Fokker-Planck equation

$$\frac{\partial \rho_t}{\partial t} = -\mathbf{grad}_W \mathcal{F}(\rho_t) = \Delta \rho_t + \nabla \cdot (\rho_t \nabla f).$$

3) In general, consider three basic kinds of energies, $\mathcal{F} = \mathcal{U} + \mathcal{V} + \mathcal{W}$, where

- internal energy $\mathcal{U}(\rho) = \int U(\rho(x)) dx$;
- potential energy $\mathcal{V}(\rho) = \int V(x) \rho(x) dx$;
- interaction energy $\mathcal{W}(\rho) = \frac{1}{2} \int W(x-y) \rho(x) \rho(y) dx dy$.

$$\frac{\delta \mathcal{W}}{\delta \rho} = \int W(\cdot - y) \rho(y) dy = W * \rho, \quad -\mathbf{grad}_W \mathcal{W}(\rho) = \nabla W * \rho.$$

Then the Wasserstein gradient flow is

$$\frac{\partial \rho_t}{\partial t} = -\mathbf{grad}_W \mathcal{F}(\rho_t) = \nabla \cdot [\rho_t \nabla U'(\rho_t) + \rho_t \nabla V + \rho_t (\rho_t * \nabla W)].$$

3 Geodesic convexity

From our construction, \mathcal{W}_2 is the geodesic length associated to the Riemannian structure (6). The geodesic is the McCann's displacement interpolation: given ρ_0 and ρ_1 , the geodesic joining them is

$$\rho_t = [(1-t)\text{Id} + t\nabla\varphi]\#\rho_0$$

where φ is a convex function and $\nabla\varphi$ is the optimal map in the Monge-Kantorovich problem from ρ_0 to ρ_1 , i.e., $\nabla\varphi\#\rho_0 = \rho_1$. The corresponding particles evolve as follows

$$X_t = (1-t)X_0 + tX_1.$$

Thus,

$$\left\| \frac{\partial\rho_t}{\partial t} \right\|_{\rho(t)}^2 = \mathbb{E}[\|\nabla\lambda(X)\|^2] = \mathbb{E}[\|c_t(X)\|^2] = \mathbb{E}[\|\dot{X}_t\|^2] = \mathbb{E}[\|X_1 - X_0\|^2] = \mathcal{W}_2^2(\rho_0, \rho_1)$$

does not depend on time. So $t \mapsto \rho_t$ is a constant speed geodesic and $\mathcal{W}_2(\rho_0, \rho_t) = t\mathcal{W}_2(\rho_0, \rho_1)$. Moreover, we have

$$c_0(x) = \nabla\varphi(x) - x, \quad c_t = c_0 \circ T_t^{-1},$$

where $T_t = (1-t)\text{Id} + t\nabla\varphi$. This can be understood from both the Eulerian and Lagrangian perspectives: $c_t(x)$ is Eulerian, can be visualized by sitting on the bank of a river and watching the water pass the fixed location; $[c_0 \circ T_t^{-1}](x)$ is Lagrangian, can be visualized as sitting in a boat and drifting down a river. The velocity at x at t is the same as that of the particle starting from some $x(0) = X_0$ passing $x(t) = x$.

Over a Riemannian manifold \mathcal{M} , the correct way to define convexity is as follows. Let \mathcal{M} be a Riemannian manifold and let \mathcal{F} be a smooth functional over \mathcal{M} . For $\alpha \in \mathbb{R}$, we say that \mathcal{F} is α -geodesically convex if one of the following equivalent conditions hold:

1. For all geodesics $(p_t)_{t \in [0,1]}$ and $t \in [0, 1]$,

$$\mathcal{F}(p_t) \leq (1-t)\mathcal{F}(p_0) + t\mathcal{F}(p_1) - \frac{\alpha t(1-t)}{2} d(p_0, p_1)^2,$$

where $d(\cdot, \cdot)$ is the induced Riemannian distance (2).

2. For all $p, q \in \mathcal{M}$,

$$\mathcal{F}(q) \geq \mathcal{F}(p) + \langle \text{grad}\mathcal{F}(p), \log_p(q) \rangle_p + \frac{\alpha}{2} d(p, q)^2. \quad (10)$$

Here, grad denotes the Riemannian gradient and $\log_p(q)$ denotes the logarithmic map $\log_p : \mathcal{M} : T_p\mathcal{M}$ taking q to the tangent vector $v \in T_p\mathcal{M}$ along which the endpoint at time 1 of the constant-speed geodesic starting from p is q . In particular, we have

$$\log_{\rho_0}(\rho_1) = -\nabla \cdot (\rho_0(\nabla\varphi - \text{Id})).$$

3. For all constant-speed geodesics $(p_t)_{t \in [0,1]}$ with tangent vector $v_0 \in T_{p_0}\mathcal{M}$ at time 0, it holds that

$$\partial_t^2 \mathcal{F}(p_t)|_{t=0} = \langle \text{Hess}\mathcal{F}(p_0)v_0, v_0 \rangle_{p_0} \geq \alpha \|v_0\|_{p_0}^2. \quad (11)$$

4 Transport inequalities

Theorem 1. $\nu \propto \exp(-f)$ and f is α -strongly convex, then $\text{KL}(\cdot|\nu)$ is α -geodesically convex along Wasserstein geodesics.

Proof. This is a well-known result. See for example page 48 of Sinho's book for the proof. \square

Theorem 2. $\nu \propto \exp(-f)$ and f is α -strongly convex, then for all $\rho_0, \rho_1 \in \mathcal{P}_2(\mathbb{R}^d)$,

$$\text{KL}(\rho_1 \parallel \nu) \geq \text{KL}(\rho_0 \parallel \nu) + \mathbb{E}_{\rho_0} \left[\left\langle \nabla \log \frac{\rho_0(x)}{\nu(x)}, \nabla \varphi(x) - x \right\rangle \right] + \frac{\alpha}{2} \mathcal{W}_2^2(\rho_0, \rho_1) \quad (12)$$

where $\nabla \varphi \# \rho_0 = \rho_1$.

Proof. It follows from Theorem 1 the definition of geodesic convexity (10) that

$$\text{KL}(\rho_1 \parallel \nu) \geq \text{KL}(\rho_0 \parallel \nu) + \langle \text{grad}_W \text{KL}(\rho_0 \parallel \nu), \log_{\rho_0}(\rho_1) \rangle_{\rho_0} + \frac{\alpha}{2} \mathcal{W}_2^2(\rho_0, \rho_1),$$

which can be reformulated to (12). \square

Taking $\rho_0 = \rho$ and $\rho_1 = \nu$ in (12), we have the *HWI inequality*

$$\text{KL}(\rho \parallel \nu) \leq \mathcal{W}_2(\rho, \nu) \sqrt{\text{FI}(\rho \parallel \nu)} - \frac{\alpha}{2} \mathcal{W}_2^2(\rho, \nu).$$

Clearly, the above inequality and the Cauchy-Schwarz inequality imply that

$$\text{KL}(\rho \parallel \nu) \leq \mathcal{W}_2(\rho, \nu) \sqrt{\text{FI}(\rho \parallel \nu)} - \frac{\alpha}{2} \mathcal{W}_2^2(\rho, \nu) \leq \frac{1}{2\alpha} \text{FI}(\rho|\nu).$$

Hence, the HWI inequality implies LSI, so we prove that if ν is α -SLC then ν satisfies α -LSI. Actually, we have a more general inequality than HWI, i.e.,

$$\text{KL}(\rho_0 \parallel \nu) \leq \text{KL}(\rho_1 \parallel \nu) + \mathcal{W}_2(\rho_0, \rho_1) \sqrt{\text{FI}(\rho_0|\nu)} - \frac{\alpha}{2} \mathcal{W}_2^2(\rho_0, \rho_1).$$

Taking $\rho_1 = \nu$, we recover HWI. Taking $\rho_0 = \nu$, we recover the Talagrand inequality (α -TI)

$$\mathcal{W}_2^2(\rho, \nu) \leq \frac{2}{\alpha} \text{KL}(\rho \parallel \nu).$$

If ν is log-concave and satisfies λ -TI, then ν also satisfies $(\lambda/4)$ -LSI. Indeed, it follows from HWI with $\alpha = 0$

$$\text{KL}(\rho \parallel \nu) \leq \mathcal{W}_2(\rho, \nu) \sqrt{\text{FI}(\rho \parallel \nu)} \leq \sqrt{\frac{2}{\lambda} \text{KL}(\rho \parallel \nu)} \sqrt{\text{FI}(\rho \parallel \nu)},$$

so

$$\text{KL}(\rho \parallel \nu) \leq \frac{2}{\lambda} \text{FI}(\rho \parallel \nu).$$

Optimization inequalities correspondence

Consider optimization in a smooth Riemannian manifold with an objective function $f : \mathcal{M} \rightarrow \mathbb{R}$, we have the following conditions guaranteeing exponential convergence:

- (1) f is α -strongly convex: $\forall x \in \mathcal{M}$, $\text{Hess}_x f \succeq \alpha I$, i.e., $(\text{Hess}_x f)(v, v) \geq \alpha \|v\|_x^2$ for all $v \in T_x \mathcal{M}$;
- (2) f is α -gradient dominant: $\forall x \in \mathcal{M}$, $\|\text{grad}_x f\|_x^2 \geq 2\alpha(f(x) - \min f)$;
- (3) f is α -sufficient growth: $\forall x \in \mathcal{M}$, $f(x) - \min f \geq \frac{\alpha}{2} d(x, x_*)^2$.

The relationship among the three inequalities is (1) \implies (2) \implies (3).

In the setting of Wasserstein space $(\mathcal{M}, d) = (\mathcal{P}_2(\mathbb{R}^d), \mathcal{W}_2)$ and the objective functional being the KL divergence $\text{KL}(\cdot \parallel \nu)$, we have corresponding functional/transport inequalities

- (a) ν is α -strongly logconcave, i.e., (11);
- (b) ν satisfies α -LSI: $\forall \rho \in \mathcal{P}_2(\mathbb{R}^d)$, $\text{FI}(\rho \parallel \nu) \geq 2\alpha \text{KL}(\rho \parallel \nu)$;
- (c) ν satisfies α -TI: $\forall \rho \in \mathcal{P}_2(\mathbb{R}^d)$, $\text{KL}(\rho \parallel \nu) \geq \frac{\alpha}{2} \mathcal{W}_2^2(\rho, \nu)$.

The same relationship also holds for the three transport inequalities: (a) \implies (b) \implies (c).

5 New convergence analysis

Now we have develop enough mathematics for Wasserstein space. We would like to study the convergence along Langevin dynamics from the Wasserstein gradient flow perspective. The following result is the same as Lemma 3 and Theorem 3 of Lecture 8, while we present a simpler derivation of the de Bruijn's identity.

Theorem 3. *Along the Langevin dynamics, we have the de Bruijn's identity*

$$\frac{d}{dt} \text{KL}(\rho_t \parallel \nu) = -\text{FI}(\rho_t \parallel \nu).$$

If ν satisfies α -LSI, we have

$$\text{KL}(\rho_t \parallel \nu) \leq e^{-2\alpha t} \text{KL}(\rho_0 \parallel \nu).$$

Proof. We know that the Fokker-Planck equation is the continuity equation of the Langevin dynamic. Moreover, we have shown that the Fokker-Planck equation is also the Wasserstein gradient flow of KL divergence

$$\frac{\partial \rho_t}{\partial t} = \Delta \rho_t + \nabla \cdot (\rho_t \nabla f) = -\text{grad}_W \text{KL}(\rho_t).$$

Hence, we have

$$\begin{aligned} \frac{d}{dt} \text{KL}(\rho_t \parallel \nu) &= \left\langle \text{grad}_W \text{KL}(\rho_t \parallel \nu), \frac{\partial \rho_t}{\partial t} \right\rangle_{\rho_t} = - \|\text{grad}_W \text{KL}(\rho_t \parallel \nu)\|_{\rho_t}^2 \\ &= -\mathbb{E}_{\rho_t} \left[\left\| \nabla \log \frac{\rho_t}{\nu} \right\|^2 \right] = -\text{FI}(\rho_t \parallel \nu). \end{aligned}$$

The rest of the proof follows that of Theorem 3 of Lecture 8. \square

Lemma 3. *Given $\nu \in \mathcal{P}_2(\mathbb{R}^d)$, the Wasserstein gradient of $\rho \mapsto \mathcal{W}_2^2(\rho, \nu)$ is $-2 \log_\rho(\nu)$. In general, on a Riemannian manifold, the gradient of $d(\cdot, q)^2$ at p is $-2 \log_p(q)$.*

Theorem 4. *If ν is log-concave, then along the Langevin dynamics, we have*

$$\text{KL}(\rho_t \parallel \nu) \leq \frac{\mathcal{W}_2^2(\rho_0, \nu)}{2t}. \quad (13)$$

Proof. Recall from Theorem 2 that

$$\text{KL}(\rho_1 \parallel \nu) \geq \text{KL}(\rho_0 \parallel \nu) + \left\langle \text{grad}_W \text{KL}(\rho_0 \parallel \nu), \log_{\rho_0}(\rho_1) \right\rangle_{\rho_0} + \frac{\alpha}{2} \mathcal{W}_2^2(\rho_0, \rho_1).$$

Taking $\rho_0 = \rho_t$, $\rho_1 = \nu$, and $\alpha = 0$, we have

$$\left\langle \text{grad}_W \text{KL}(\rho_t \parallel \nu), \log_{\rho_t}(\nu) \right\rangle_{\rho_t} \leq -\text{KL}(\rho_t \parallel \nu). \quad (14)$$

Consider a Lyapunov functional $\mathcal{L}_t := t \text{KL}(\rho_t \parallel \nu) + \frac{1}{2} \mathcal{W}_2^2(\rho_t, \nu)$. Taking time derivative gives

$$\begin{aligned} \frac{\partial}{\partial t} \mathcal{L}_t &= \text{KL}(\rho_t \parallel \nu) - t \text{FI}(\rho_t \parallel \nu) + \left\langle \frac{\partial \rho_t}{\partial t}, -\log_{\rho_t}(\nu) \right\rangle_{\rho_t} \\ &= \text{KL}(\rho_t \parallel \nu) - t \text{FI}(\rho_t \parallel \nu) + \left\langle \text{grad}_W \text{KL}(\rho_t \parallel \nu), \log_{\rho_t}(\nu) \right\rangle_{\rho_t} \\ &\leq -t \text{FI}(\rho_t \parallel \nu) \leq 0 \end{aligned}$$

where the first inequality is due to (14). Therefore, $\mathcal{L}_t \leq \mathcal{L}_0$ and (13) holds. \square

We provide a convergence guarantee of the proximal sampling algorithm in Lecture 9 for the log-concave case, which closely resembles the convergence of proximal point method in optimization for convex functions.

Theorem 5. *If ν is log-concave, then ρ_k^X of the proximal sampling algorithm satisfies*

$$\text{KL}(\rho_k^X \parallel \nu) \leq \frac{\mathcal{W}_2^2(\rho_0^X, \nu)^2}{k\eta}.$$

Proof. Since ν_t is log-concave (log-concavity is preserved by convolution), we have HWI inequality

$$\text{KL}(\rho_t \parallel \nu_t) \leq \mathcal{W}_2(\rho_t, \nu_t) \sqrt{\text{FI}(\rho_t \parallel \nu_t)}.$$

It follows from Lemma 3 of Lecture 9 and the above HWI inequality that

$$\frac{d}{dt} \text{KL}(\rho_t \parallel \nu_t) = -\frac{1}{2} \text{FI}(\rho_t \parallel \nu_t) \leq -\frac{1}{2} \frac{\text{KL}(\rho_t \parallel \nu_t)^2}{\mathcal{W}_2^2(\rho_t, \nu_t)} \leq -\frac{1}{2} \frac{\text{KL}(\rho_t \parallel \nu_t)^2}{\mathcal{W}_2^2(\rho_0, \nu_0)} = -\frac{1}{2} \frac{\text{KL}(\rho_t \parallel \nu_t)^2}{\mathcal{W}_2^2(\rho_k^X, \pi^X)},$$

where the last inequality is due to the fact that \mathcal{W}_2 contraction along heat flow (proof by coupling). Solving this differential inequality yields

$$\frac{1}{\text{KL}(\rho_k^Y \parallel \pi^Y)} = \frac{1}{\text{KL}(\rho_\eta \parallel \nu_\eta)} \geq \frac{1}{\text{KL}(\rho_k^X \parallel \pi^X)} + \frac{\eta}{2\mathcal{W}_2^2(\rho_k^X, \pi^X)}. \quad (15)$$

For the backward process, we similarly have

$$\frac{d}{dt} \text{KL}(\rho_t^- \parallel \nu_t^-) = -\frac{1}{2} \text{FI}(\rho_t^- \parallel \nu_t^-) \leq -\frac{1}{2} \frac{\text{KL}(\rho_t^- \parallel \nu_t^-)^2}{\mathcal{W}_2^2(\rho_t^-, \nu_t^-)}.$$

Recall from Lecture 9 that the backward channel can be modeled by the following SDE

$$dY_t = \nabla \log \nu_{\eta-t}(Y_t) dt + dW_t. \quad (16)$$

Since $\log \nu_{\eta-t}$ is log-concave, by a coupling argument, we can show $t \mapsto \mathcal{W}_2^2(\rho_t^-, \nu_t^-)$ is decreasing. Thus,

$$\mathcal{W}_2^2(\rho_t^-, \nu_t^-) \leq \mathcal{W}_2^2(\rho_0^-, \nu_0^-) = \mathcal{W}_2^2(\rho_k^Y, \pi^Y) \leq \mathcal{W}_2^2(\rho_k^X, \pi^X).$$

Therefore, we obtain

$$\frac{1}{\text{KL}(\rho_{k+1}^X \parallel \pi^X)} = \frac{1}{\text{KL}(\rho_\eta^- \parallel \nu_\eta^-)} \geq \frac{1}{\text{KL}(\rho_k^Y \parallel \pi^Y)} + \frac{\eta}{2\mathcal{W}_2^2(\rho_k^X, \pi^X)}.$$

Combining the above inequality and (15), we have

$$\frac{1}{\text{KL}(\rho_{k+1}^X \parallel \pi^X)} \geq \frac{1}{\text{KL}(\rho_k^X \parallel \pi^X)} + \frac{\eta}{\mathcal{W}_2^2(\rho_k^X, \pi^X)} \geq \frac{1}{\text{KL}(\rho_k^X \parallel \pi^X)} + \frac{\eta}{\mathcal{W}_2^2(\rho_0^X, \pi^X)},$$

i.e.,

$$\frac{1}{\text{KL}(\rho_{k+1}^X \parallel \nu)} \geq \frac{1}{\text{KL}(\rho_k^X \parallel \nu)} + \frac{\eta}{\mathcal{W}_2^2(\rho_0^X, \nu)}.$$

Summing the above inequality over iterations, we arrive at

$$\frac{1}{\text{KL}(\rho_k^X \parallel \nu)} \geq \frac{1}{\text{KL}(\rho_0^X \parallel \nu)} + \frac{k\eta}{\mathcal{W}_2^2(\rho_0^X, \nu)} \geq \frac{k\eta}{\mathcal{W}_2^2(\rho_0^X, \nu)}.$$

□

6 Revisiting proximal sampling

Recall that sampling from $\nu \propto \exp(-f)$ can be understood as minimizing $\text{KL}(\cdot \parallel \nu)$ over $\mathbb{R}_2(\mathbb{R}^d)$. Because the Fokker-Planck equation is the Wasserstein gradient flow of KL divergence, in the language of numerical methods, the Langevin dynamics (and hence LMC) can be viewed as an explicit scheme. As we know explicit schemes are less stable than implicit schemes, and an example in optimization is gradient descent versus proximal point method. In sampling, the implicit scheme is known as the JKO scheme, which repeatedly invokes the proximal operator over the Wasserstein space

$$\rho_{k+1} = \text{prox}_{\eta\mathcal{F}}(\rho_k) = \underset{\rho \in \mathcal{P}_2(\mathbb{R}^d)}{\text{argmin}} \left\{ \mathcal{F}(\rho) + \frac{1}{2\eta} \mathcal{W}_2^2(\rho_k, \rho) \right\}.$$

Hence, the JKO scheme is a Wasserstein analogue of the proximal point method.

We revisit the proximal sampling algorithm in Lecture 9 from this variational perspective. The first result shows that RGO is a proximal operator on the Wasserstein space.

Lemma 4. *Recall RGO is a sampling oracle for given $y \in \mathbb{R}^d$ that returns $x \sim \pi^{X|Y}(x|y)$ where*

$$\pi^{X|Y}(x | y) \propto \exp \left(-f(x) - \frac{1}{2\eta} \|x - y\|^2 \right).$$

Then, RGO is a proximal operator of KL divergence, i.e.,

$$\pi^{X|Y=y} = \text{prox}_{\eta\text{KL}(\cdot \parallel \pi^X)}(\delta_y).$$

Proof. Since $\pi^X(x) = \nu(x) \propto \exp(-f(x))$, we know

$$\pi^{X|Y}(x | y) \propto \exp \left(-\frac{1}{2\eta} \|x - y\|^2 \right) \pi^X(x).$$

Also, note that

$$\begin{aligned} \text{KL}(\rho^X \parallel \pi^{X|Y}) &= \int \rho^X \log \frac{\rho^X}{\pi^{X|Y}} = \int \rho^X \left(\log \frac{\rho^X}{\pi^X} + \frac{1}{2\eta} \|x - y\|^2 \right) + C(y) \\ &= \text{KL}(\rho^X \parallel \pi^X) + \int \rho^X \frac{1}{2\eta} \|x - y\|^2 + C(y), \end{aligned}$$

where $C(y)$ is a function of y . Then, RGO can be expressed as

$$\begin{aligned} \pi^{X|Y=y} &= \underset{\rho^X \in \mathcal{P}_2(\mathbb{R}^d)}{\text{argmin}} \text{KL}(\rho^X \parallel \pi^{X|Y}) = \underset{\rho^X \in \mathcal{P}_2(\mathbb{R}^d)}{\text{argmin}} \left\{ \text{KL}(\rho^X \parallel \pi^X) + \frac{1}{2\eta} \int \|x - y\|^2 d\rho^X(x) \right\} \\ &= \underset{\rho^X \in \mathcal{P}_2(\mathbb{R}^d)}{\text{argmin}} \left\{ \text{KL}(\rho^X \parallel \pi^X) + \frac{1}{2\eta} \mathcal{W}_2^2(\rho^X, \delta_y) \right\} = \text{prox}_{\eta\text{KL}(\cdot \parallel \pi^X)}(\delta_y). \end{aligned}$$

□

A variant of the JKO scheme introduces an extra entropic regularization term

$$\rho_{k+1} = \operatorname{argmin}_{\rho \in \mathcal{P}_2(\mathbb{R}^d)} \left\{ \mathcal{F}(\rho) + \frac{1}{2\eta} \mathcal{W}_{2,\varepsilon}^2(\rho_k, \rho) \right\}$$

where $\mathcal{W}_{2,\varepsilon}$ is the entropy-regularized Wasserstein-2 distance defined as

$$\mathcal{W}_{2,\varepsilon}^2(\mu, \nu) := \min_{\gamma \in \Gamma(\mu, \nu)} \left\{ \int \|x - y\|^2 \gamma(x, y) dx dy - \varepsilon \mathbf{H}(\gamma) \right\} \quad (17)$$

and $\mathbf{H}(\gamma) = - \int \gamma \log \gamma$ denotes the entropy.

The next result shows that the proximal sampling algorithm can be viewed as the entropy-regularized JKO scheme.

Theorem 6. *Let $\rho_k^X, \rho_k^Y, \rho_{k+1}^X$ be the distributions of x_k, y_k, x_{k+1} , respectively, in one iteration of the proximal sampling algorithm. Then, they follow the entropy-regularized JKO scheme*

$$\rho_k^Y = \operatorname{argmin}_{\rho \in \mathcal{P}_2(\mathbb{R}^d)} \frac{1}{2\eta} \mathcal{W}_{2,2\eta}^2(\rho_k^X, \rho), \quad (18)$$

$$\rho_{k+1}^X = \operatorname{argmin}_{\rho \in \mathcal{P}_2(\mathbb{R}^d)} \left\{ \int f \rho + \frac{1}{2\eta} \mathcal{W}_{2,2\eta}^2(\rho_k^Y, \rho) \right\}. \quad (19)$$

Proof. Plugging (17) into the right-hand side of (18) gives a solution γ^Y with γ being the solution to

$$\min_{\substack{\gamma \in \mathcal{P}_2(\mathbb{R}^d \times \mathbb{R}^d) \\ \gamma^X = \rho_k^X}} \left\{ \int \frac{1}{2\eta} \|x - y\|^2 \gamma(x, y) dx dy - \mathbf{H}(\gamma) \right\}.$$

Consider the Lagrangian function

$$\begin{aligned} & \int \frac{1}{2\eta} \|x - y\|^2 \gamma(x, y) dx dy - \mathbf{H}(\gamma) + \int \lambda(x) [\gamma^X(x) - \rho_k^X(x)] dx \\ &= \int \frac{1}{2\eta} \|x - y\|^2 \gamma(x, y) dx dy - \mathbf{H}(\gamma) + \int \lambda(x) [\gamma(x, y) - \rho_k^X(x)] dx dy \end{aligned}$$

Taking the first variation yields

$$\frac{1}{2\eta} \|x - y\|^2 + 1 + \log \gamma + \lambda(x) = 0.$$

So

$$\gamma(x, y) \propto \exp \left(-\frac{1}{2\eta} \|x - y\|^2 - \lambda(x) \right),$$

and the X-marginal is $\gamma^X \propto \exp(-\lambda(x))$. On the other hand, from the constraint, we know $\gamma^X = \rho_k^X$. So

$$\gamma(x, y) \propto \exp \left(-\frac{1}{2\eta} \|x - y\|^2 \right) \rho_k^X(x)$$

and $\gamma^Y = \rho_k^X * \mathcal{N}(0, \eta I) = \rho_k^Y$. Thus, we verify (18).

Similarly, plugging (17) into the right-hand side of (19) gives a solution γ^X with γ being the solution to

$$\operatorname{argmin}_{\substack{\gamma \in \mathcal{P}_2(\mathbb{R}^d \times \mathbb{R}^d) \\ \gamma^Y = \rho_k^Y}} \left\{ \int \left[f(x) + \frac{1}{2\eta} \|x - y\|^2 \right] \gamma(x, y) dx dy - \mathbf{H}(\gamma) \right\}.$$

Clearly, $\gamma(x, y) \propto \rho_k^Y(y) \exp\left(-f(x) - \frac{1}{2\eta} \|x - y\|^2\right)$. Hence, the X-marginal γ^X has the same distribution as ρ_{k+1}^X . We thus complete the proof. \square

The above interpretation of the proximal sampling algorithm provides some insights on its connections to optimization. Define a more general proximal sampling algorithm with a different level of entropy regularization

$$\rho_k^Y = \operatorname{argmin}_{\rho \in \mathcal{P}_2(\mathbb{R}^d)} \frac{1}{2\eta} \mathcal{W}_{2, 2\eta\varepsilon}^2(\rho_k^X, \rho), \quad (20)$$

$$\rho_{k+1}^X = \operatorname{argmin}_{\rho \in \mathcal{P}_2(\mathbb{R}^d)} \left\{ \int f \rho + \frac{1}{2\eta} \mathcal{W}_{2, 2\eta\varepsilon}^2(\rho_k^Y, \rho) \right\}. \quad (21)$$

The next result reveals the well-known fact that optimization is the limit of sampling on the particular example of proximal sampling and optimization.

Theorem 7. *As $\varepsilon \rightarrow 0$, (20) and (21) reduce to the proximal point method in optimization.*

Proof. Following a similar argument as in the proof of Theorem 6, we can show that (20)-(21) correspond to the following proximal sampling

$$\begin{aligned} y_k &\sim \pi_\varepsilon^{Y|X=x_k} = \mathcal{N}(x_k, \varepsilon\eta I), \\ x_{k+1} &\sim \pi_\varepsilon^{X|Y=y_k} \propto \exp\left[-\frac{1}{\varepsilon} \left(f(x) + \frac{1}{2\eta} \|x - y_k\|^2 \right)\right]. \end{aligned}$$

As $\varepsilon \rightarrow 0$, the above steps becomes to

$$\begin{aligned} y_k &= x_k, \\ x_{k+1} &= \operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ f(x) + \frac{1}{2\eta} \|x - y_k\|^2 \right\}. \end{aligned}$$

Putting them together, we have $x_{k+1} = \operatorname{prox}_{\eta f}(x_k)$, i.e., the proximal point method. \square

We finally note that the stationary distribution of the new proximal sampling algorithm (20)-(21) is $\pi_\varepsilon^X \propto \exp(-f/\varepsilon)$, which converges (as $\varepsilon \rightarrow 0$) to a Dirac distribution concentrating on the minimizer of f (or a uniform distribution over the minimizer set of f).