

Supplementary Materials for “Learning Gaussian Mixtures Using the Wasserstein-Fisher-Rao Gradient Flow”

April 21, 2023

A Preliminaries

In the main text, we focused on the Gaussian mixture model where $\phi(x) = (2\pi)^{-d/2} \exp(-\|x\|_2^2/2)$. In fact, the algorithms and theorems in the current paper are also valid for a more general class of probability density ϕ . In the appendices, we only assume that ϕ satisfies the following regularity assumption.

Assumption 1 (Regularity). *Assume that the density $\phi(x) > 0$ for any $x \in \mathbb{R}^d$. Furthermore, $\phi \in C^{\max\{d, 2\}}(\mathbb{R}^d)$, $\lim_{\|x\|_2 \rightarrow \infty} \phi(x) = 0$, $\sup_{x \in \mathbb{R}^d} \|\nabla \phi(x)\|_2 < \infty$ and $\sup_{x \in \mathbb{R}^d} \|\nabla^2 \phi(x)\|_2 < \infty$.*

It is clear that the Gaussian kernel $\phi(x) = (2\pi)^{-d/2} \exp(-\|x\|_2^2/2)$ satisfies Assumption 1. The following lemma shows that NPMLE is compactly supported. The proof can be found in Appendix B.2.

Lemma 1 (Compact support of NPMLE). *Let Assumption 1 hold and $\hat{\rho}$ be any optimal solution to (1.1). Define $R_1 = \inf\{r \geq 0 : \bar{\phi}(r) \leq \underline{\phi}[\text{diam}(\Omega)]/2\}$ and*

$$R = \inf \left\{ r \geq 0 : \bar{\phi}(r) \leq \frac{\bar{\phi}(R_1) \underline{\phi}(R_1 + \text{diam}(\Omega))}{8\bar{\phi}(0)} \right\}.$$

Then, we have $\text{supp}(\hat{\rho}) \subseteq \Omega_R$.

B Proof of structural results for NPMLE

In this section, we prove the two structural results for NPMLE, namely Theorem 1 and Lemma 1 under Assumption 1.

B.1 Proof of Theorem 1

Part 1: existence of NPMLE. Note that the loss function ℓ_N is lower bounded

$$\ell_N(\rho) = -\frac{1}{N} \sum_{i=1}^N \log[(\rho * \phi)(X_i)] \geq -\log \|\phi\|_\infty, \quad (\text{B.1})$$

where the last inequality follows from $\rho * \phi(x) = \int_{\mathbb{R}^d} \phi(y - x) \rho(dy) \leq \|\phi\|_\infty$ for any $x \in \mathbb{R}^d$. Therefore there exists a sequence of probability distribution $\{\rho_n\}$ such that

$$\ell_N(\rho_n) \leq \inf_{\rho \in \mathcal{P}(\mathbb{R}^d)} \ell_N(\rho) + \frac{1}{n}. \quad (\text{B.2})$$

Now we argue that $\{\rho_n\}$ is tight. To that end, we show that there exists $r > 0$ such that for any $\varepsilon > 0$, it holds $\rho_n(\Omega_r) \geq 1 - \varepsilon$ for n large enough.

For any n and $r > 0$, define

$$\rho_{n,r} := \rho_n(\Omega_r) \cdot \rho_n|_{\Omega_r} + \rho_n(\Omega_r^c) \cdot \text{Unif}(\Omega),$$

where $\rho_n|_{\Omega_r}(\cdot) = \rho_n(\cdot|\Omega_r)$ is the conditional distribution of ρ_n given Ω_r , and $\text{Unif}(\Omega)$ is the uniform distribution on Ω . We have

$$\ell_N(\rho_n) - \ell_N(\rho_{n,r}) = -\frac{1}{N} \sum_{i=1}^N \log[(\rho_n * \phi)(X_i)] + \frac{1}{N} \sum_{i=1}^N \log[(\rho_{n,r} * \phi)(X_i)] = \frac{1}{N} \sum_{i=1}^N \log \left[\frac{(\rho_{n,r} * \phi)(X_i)}{(\rho_n * \phi)(X_i)} \right].$$

Note that for each $i \in [N]$

$$\begin{aligned} \log \left[\frac{(\rho_{n,r} * \phi)(X_i)}{(\rho_n * \phi)(X_i)} \right] &= \log \left[\frac{\int_{\Omega_r} \phi(X_i - y) \rho_n(dy) + \rho_n(\Omega_r^c) \int_{\Omega} \phi(X_i - y) dy}{\int_{\Omega_r} \phi(X_i - y) \rho_n(dy) + \int_{\Omega_r^c} \phi(X_i - y) \rho_n(dy)} \right] \\ &= \log \left[1 + \frac{\rho_n(\Omega_r^c) \int_{\Omega} \phi(X_i - y) dy - \int_{\Omega_r^c} \phi(X_i - y) \rho_n(dy)}{(\rho_n * \phi)(X_i)} \right] \\ &\geq \log \left[1 + \frac{\rho_n(\Omega_r^c) [\phi(\text{diam}(\Omega)) - \bar{\phi}(r)]}{\|\phi\|_{\infty}} \right]. \end{aligned}$$

We can choose $r > 0$ to be sufficiently large so that $\bar{\phi}(r) \leq \phi(\text{diam}(\Omega))/2$, and therefore for each $i \in [N]$

$$\ell_N(\rho_n) - \ell_N(\rho_{n,r}) \geq \log \left[1 + \frac{\rho_n(\Omega_r^c) \phi(\text{diam}(\Omega))}{2 \|\phi\|_{\infty}} \right].$$

In view of (B.2), we know that

$$\ell_N(\rho_n) - \ell_N(\rho_{n,r}) \leq \frac{1}{n}.$$

Taking the above two inequalities collectively give

$$\rho_n(\Omega_r^c) \leq \frac{2 \|\phi\|_{\infty} [\exp(\frac{1}{n}) - 1]}{\phi(\text{diam}(\Omega))}.$$

Therefore we have

$$\rho_n(\Omega_r^c) \leq \varepsilon \quad \text{as long as} \quad n \geq n_{\varepsilon} := \left\lceil 1 / \log \left[1 + \frac{\varepsilon \phi(\text{diam}(\Omega))}{2 \|\phi\|_{\infty}} \right] \right\rceil,$$

which implies that $\{\rho_n\}$ is tight. We conclude using Prokhorov's theorem: there exists a subsequence $\{\rho_{n_k}\}$ and $\hat{\rho} \in \mathcal{P}(\mathbb{R}^d)$ such that ρ_{n_k} converges weakly to $\hat{\rho}$ which must be a minimizer of (1.1).

Part 2: optimality condition. First of all, it is straightforward to check that for any $\rho \in \mathcal{M}(\mathbb{R}^d)$

$$\int_{\mathbb{R}^d} \delta \ell_N(\rho)(x) \rho(dx) = -\frac{1}{N} \sum_{i=1}^N \frac{\int \phi(x - X_i) \rho(dx)}{(\rho * \phi)(X_i)} = -\frac{1}{N} \sum_{i=1}^N \frac{(\rho * \phi)(X_i)}{(\rho * \phi)(X_i)} = -1.$$

If $\hat{\rho} \in \mathcal{M}(\mathbb{R}^d)$ is the optimal solution to (1.1), then for any $x \in \mathbb{R}^d$ and any $\varepsilon \in [0, 1]$ we have

$$\ell_N(\hat{\rho}) \leq \ell_N((1 - \varepsilon)\hat{\rho} + \varepsilon\delta_x) = \ell_N(\rho + \varepsilon(\delta_x - \hat{\rho})).$$

As a result we have

$$\delta \ell_N(\hat{\rho})(x) + 1 = \int_{\mathbb{R}^d} \delta \ell_N(\hat{\rho}) d(\delta_x - \hat{\rho}) = \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} [\ell_N(\hat{\rho} + \varepsilon(\delta_x - \hat{\rho})) - \ell_N(\hat{\rho})] \geq 0.$$

Since x is arbitrary, this implies that $\delta \ell_N(\hat{\rho})(x) \geq -1$ for any $x \in \mathbb{R}^d$. Combine this with $\int_{\mathbb{R}^d} \delta \ell_N(\hat{\rho}) d\hat{\rho} = -1$ readily gives $\delta \ell_N(\hat{\rho})(x) = -1$ for $\hat{\rho}$ -a.e. x .

Conversely, if $\hat{\rho} \in \mathcal{M}(\mathbb{R}^d)$ satisfies $\delta \ell_N(\hat{\rho})(x) \geq -1$ for all $x \in \mathbb{R}^d$, then for any $\rho \in \mathcal{M}(\mathbb{R}^d)$, it holds

$$0 \leq \int_{\mathbb{R}^d} \delta \ell_N(\hat{\rho}) d\rho + 1 = \int_{\mathbb{R}^d} \delta \ell_N(\hat{\rho}) d(\rho - \hat{\rho}) = \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} [\ell_N(\hat{\rho} + \varepsilon(\rho - \hat{\rho})) - \ell_N(\hat{\rho})] \leq \ell_N(\rho) - \ell_N(\hat{\rho})$$

where the last inequality follows from convexity of the functional $\rho \mapsto \ell_N(\rho)$. The above display yields that $\hat{\rho}$ is a global minimizer of (1.1).

B.2 Proof of Lemma 1

By Theorem 1, $\delta\ell_N(\hat{\rho}) = -1$ over $\text{supp}(\hat{\rho})$. We will show that $|\delta\ell_N(\hat{\rho})(y)| < 1/2$ when y is too far away from Ω , and then conclude that $\text{supp}(\hat{\rho})$ must stay close to Ω . To that end, we present some useful estimates in the following lemma that will also be useful later.

Lemma 2. *Let Assumption 1 hold. For any $\rho \in \mathcal{P}(\mathbb{R}^d)$, $x \in \Omega$ and $r \geq 0$, we have*

$$\rho(\Omega_r)\phi(r + \text{diam}(\Omega)) \leq (\rho * \phi)(x) \leq \rho(\Omega_r^c)\bar{\phi}(r) + \rho(\Omega_r)\bar{\phi}(0) \leq \bar{\phi}(r) + \rho(\Omega_r)\bar{\phi}(0).$$

As a result,

$$-\log(\bar{\phi}(r) + \rho(\Omega_r)\bar{\phi}(0)) \leq \ell_N(\rho) \leq -\log(\rho(\Omega_r)\phi(r + \text{diam}(\Omega))).$$

Take any $R \geq 0$ such that $\bar{\phi}(R) \leq e^{-\ell_N(\rho)}/2$. For any $\mu \in \mathcal{P}(\mathbb{R}^d)$ with $\ell_N(\mu) \leq \ell_N(\rho)$, we have

$$\begin{aligned} \mu(\Omega_R) &\geq e^{-\ell_N(\rho)}/[2\bar{\phi}(0)], \\ \inf_{x \in \Omega} (\mu * \phi)(x) &\geq e^{-\ell_N(\rho)}\phi(R + \text{diam}(\Omega))/[2\bar{\phi}(0)], \\ \sup_{y \in \Omega_r^c} |\delta\ell_N(\mu)(y)| &\leq \frac{2e^{\ell_N(\rho)}\bar{\phi}(0)}{\phi(R + \text{diam}(\Omega))} \cdot \bar{\phi}(r), \quad \forall r \geq 0. \end{aligned}$$

The proof of Lemma 2 is deferred to the end of this section. We come back to proving Lemma 1. Choose any $\rho_0 \in \mathcal{P}(\mathbb{R}^d)$ supported on Ω . By Lemma 2, we have

$$\ell_N(\rho_0) \leq -\log(\phi[\text{diam}(\Omega)]), \quad \forall r \geq 0.$$

Take $R_1 = \inf\{r \geq 0 : \bar{\phi}(r) \leq \phi[\text{diam}(\Omega)]/2\}$. By the continuity of $\bar{\phi}$,

$$\bar{\phi}(R_1) \leq \phi[\text{diam}(\Omega)]/2 \leq e^{-\ell_N(\rho_0)}/2.$$

Lemma 2 implies that

$$\sup_{y \in \Omega_r^c} |\delta\ell_N(\hat{\rho})(y)| \leq \frac{2e^{\ell_N(\rho_0)}\bar{\phi}(0)}{\phi(R_1 + \text{diam}(\Omega))} \cdot \bar{\phi}(r) \leq \frac{4\bar{\phi}(0)}{\bar{\phi}(R_1)\phi(R_1 + \text{diam}(\Omega))} \cdot \bar{\phi}(r), \quad \forall r \geq 0.$$

Let

$$R = \inf \left\{ r \geq 0 : \bar{\phi}(r) \leq \frac{\bar{\phi}(R_1)\phi(R_1 + \text{diam}(\Omega))}{8\bar{\phi}(0)} \right\}.$$

The continuity of $\bar{\phi}$ leads to $\bar{\phi}(R) \leq \bar{\phi}(R_1)\phi(R_1 + \text{diam}(\Omega))/[8\bar{\phi}(0)]$ and thus $\sup_{y \in \Omega_R^c} |\delta\ell_N(\hat{\rho})(y)| \leq 1/2 < 1$. As a result, $\text{supp}(\hat{\rho}) \cap \Omega_R^c = \emptyset$ and $\text{supp}(\hat{\rho}) \subseteq \Omega_R$.

Proof of Lemma 2. Note that $\phi(x - y) \leq \bar{\phi}(0)$ for any $x, y \in \mathbb{R}^d$. If $x \in \Omega$ and $y \in \Omega_r^c$, then $\|x - y\|_2 \geq r$ and thus $\phi(x - y) \leq \bar{\phi}(r)$. Hence,

$$(\rho * \phi)(x) \leq \int_{\Omega_r} \phi(x - y)\rho(dy) + \int_{\Omega_r^c} \phi(x - y)\rho(dy) \leq \bar{\phi}(0)\rho(\Omega_r) + \bar{\phi}(r)\rho(\Omega_r^c).$$

If $x \in \Omega$ and $y \in \Omega_r$, then $\|x - y\|_2 \geq r + \text{diam}(\Omega)$ and thus $\phi(x - y) \geq \phi(r + \text{diam}(\Omega))$. Therefore,

$$(\rho * \phi)(x) \geq \int_{\Omega_r} \phi(x - y)\rho(dy) \geq \phi(r + \text{diam}(\Omega))\rho(\Omega_r).$$

The desired bounds on $\ell_N(\cdot)$ become obvious. If $\mu \in \mathcal{P}(\mathbb{R}^d)$ and $\ell_N(\mu) \leq \ell_N(\rho)$, then our estimates of ℓ implies that

$$-\log(\bar{\phi}(r) + \mu(\Omega_r)\bar{\phi}(0)) \leq \ell_N(\mu) \leq \ell_N(\rho), \quad \forall r \geq 0.$$

Hence, $\bar{\phi}(r) + \mu(\Omega_r)\bar{\phi}(0) \geq e^{-\ell_N(\rho)}$. By Assumption 1, $\lim_{r \rightarrow \infty} \bar{\phi}(r) = 0$. Take any $R \geq 0$ such that $\bar{\phi}(R) \leq e^{-\ell_N(\rho)}/2$. Then,

$$\begin{aligned}\mu(\Omega_R) &\geq [e^{-\ell_N(\rho)} - \bar{\phi}(R)]/\bar{\phi}(0) \geq e^{-\ell_N(\rho)}/[2\bar{\phi}(0)], \\ \inf_{x \in \Omega} (\mu * \phi)(x) &\geq \mu(\Omega_R)\bar{\phi}(R + \text{diam}(\Omega)) \geq e^{-\ell_N(\rho)}\bar{\phi}(R + \text{diam}(\Omega))/[2\bar{\phi}(0)].\end{aligned}$$

For any $r \geq 0$, we have $\|X - y\|_2 \geq r$ whenever $X \in \text{supp}(\nu)$ and $y \in \Omega_r^c$. Then,

$$|\delta \ell_N(\mu)(y)| = \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - y)}{(\mu * \phi)(X_i)} \leq \frac{\bar{\phi}(r)}{\inf_{x \in \Omega} (\mu * \phi)(x)}, \quad \forall r \geq 0, \quad y \in \Omega_r^c.$$

The proof is finished by combining this and the lower bound on $\inf_{x \in \Omega} (\mu * \phi)(x)$ we have established above. \square

C Derivation of gradient flows over $\mathcal{P}_2(\mathbb{R}^d)$

C.1 First variation

Recall that the population and finite-sample loss functions are

$$\ell_\infty(\rho) = \mathbb{E}_{X \sim (\rho * \phi)} \{ \log [(\rho * \phi)(X)] \}, \quad \ell_N(\rho) = -\frac{1}{N} \sum_{i=1}^N \log [(\rho * \phi)(X_i)].$$

The first variation of ℓ_N is defined as any measurable function $\delta \ell(\rho) : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfying

$$\lim_{\varepsilon \rightarrow 0} \frac{\ell(\rho + \varepsilon \mathcal{X}) - \ell(\rho)}{\varepsilon} = \int \delta \ell(\rho) d\mathcal{X}$$

for all signed measures \mathcal{X} satisfying $\int d\mathcal{X} = 0$. In particular, it is easy to see that the first variation is defined up to an additive constant.

By direct computation, we have

$$\begin{aligned}\lim_{\varepsilon \rightarrow 0} \frac{\ell_N(\rho + \varepsilon \mathcal{X}) - \ell_N(\rho)}{\varepsilon} &= -\frac{1}{N} \sum_{i=1}^N \lim_{\varepsilon \rightarrow 0} \frac{\log [(\rho + \varepsilon \mathcal{X}) * \phi](X_i) - \log [(\rho * \phi)(X_i)]}{\varepsilon} \\ &= -\frac{1}{N} \sum_{i=1}^N \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} \log \left[1 + \varepsilon \frac{(\mathcal{X} * \phi)(X_i)}{(\rho * \phi)(X_i)} \right] = -\frac{1}{N} \sum_{i=1}^N \frac{(\mathcal{X} * \phi)(X_i)}{(\rho * \phi)(X_i)} \\ &= -\frac{1}{N} \sum_{i=1}^N \int \frac{\phi(x - X_i)}{(\rho * \phi)(X_i)} \mathcal{X}(dx),\end{aligned}$$

As a result, we have

$$\delta \ell_N(\rho) : x \rightarrow -\frac{1}{N} \sum_{i=1}^N \frac{\phi(x - X_i)}{(\rho * \phi)(X_i)}.$$

Similarly, we can also compute

$$\begin{aligned}\lim_{\varepsilon \rightarrow 0} \frac{\ell_\infty(\rho + \varepsilon \mathcal{X}) - \ell_\infty(\rho)}{\varepsilon} &= \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} \left[-\int \log \left[\frac{(\rho + \varepsilon \mathcal{X}) * \phi(x)}{(\rho * \phi)(x)} \right] (\rho^* * \phi)(x) dx \right] \\ &= -\int \frac{(\mathcal{X} * \phi)(x)}{(\rho * \phi)(x)} (\rho^* * \phi)(x) dx = -\int \frac{(\rho^* * \phi)(x)}{(\rho * \phi)(x)} \left[\int \phi(x - y) \mathcal{X}(dy) \right] dx \\ &= -\int \left[\int \frac{(\rho^* * \phi)(x)}{(\rho * \phi)(x)} \phi(x - y) dx \right] \mathcal{X}(dy),\end{aligned}$$

which gives

$$\delta \ell_\infty(\rho) : x \rightarrow -\int \frac{(\rho^* * \phi)(y)}{(\rho * \phi)(y)} \phi(x - y) dy.$$

C.2 Fisher-Rao gradient flow

C.2.1 A formal derivation of gradient flow using Riemannian geometry

We first introduce a Riemannian structure over $\mathcal{P}_2(\mathbb{R}^d)$ underlying the Fisher-Rao metric. Define the tangent space at $\rho \in \mathcal{P}_2(\mathbb{R}^d)$ as

$$\text{Tan}_\rho^{\text{FR}} \mathcal{P}_2(\mathbb{R}^d) := \left\{ \zeta : \zeta = \rho \left(\alpha - \int \alpha d\rho \right) \text{ for some } \alpha \text{ satisfying } \int \alpha^2 d\rho < \infty \right\}.$$

We equip the tangent space $\text{Tan}_\rho^{\text{FR}} \mathcal{P}_2(\mathbb{R}^d)$ with the following Riemannian metric tensor $g_\rho^{\text{FR}}(\cdot, \cdot) : \text{Tan}_\rho^{\text{FR}} \mathcal{P}_2(\mathbb{R}^d) \times \text{Tan}_\rho^{\text{FR}} \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ as

$$\begin{aligned} g_\rho^{\text{FR}}(\zeta_1, \zeta_2) &:= \int \frac{\zeta_1 \cdot \zeta_2}{\rho^2} d\rho \\ &= \int_{\mathbb{R}^d} \left[\alpha_1(x) - \int_{\mathbb{R}^d} \alpha_1 d\rho \right] \left[\alpha_2(x) - \int_{\mathbb{R}^d} \alpha_2 d\rho \right] \rho(dx) \\ &= \int_{\mathbb{R}^d} \alpha_1(x) \alpha_2(x) \rho(dx) - \left(\int_{\mathbb{R}^d} \alpha_1 d\rho \right) \left(\int_{\mathbb{R}^d} \alpha_2 d\rho \right), \end{aligned}$$

for any $\zeta_1 = \rho(\alpha_1 - \int \alpha_1 d\rho)$ and $\zeta_2 = \rho(\alpha_2 - \int \alpha_2 d\rho)$. The metric induced by this Riemannian structure, namely the Fisher-Rao metric $d_{\text{FR}}(\cdot, \cdot)$, satisfies the following property:

$$d_{\text{FR}}^2(\rho_0, \rho_1) = \inf \left\{ \int_0^1 \int \left[\left(\alpha_t - \int \alpha_t d\rho_t \right)^2 \right] d\rho_t dt : (\rho_t, \alpha_t)_{t \in [0,1]} \text{ solves } \partial_t \rho_t = \rho_t \alpha_t \right\}.$$

Then we follow [Gallouët and Monsaingeon \(2017\)](#); [Lu et al. \(2019\)](#) to derive the Fisher-Rao gradient flow with respect to the functional ℓ_N . Let $(\rho_t)_{t \geq 0}$ be a C^1 curve satisfying $\rho_0 = \rho$ with initial velocity

$$\partial_t \rho_t|_{t=0} = \zeta = \rho \left(\alpha - \int \alpha d\rho \right).$$

The Fisher-Rao gradient of ℓ_N at ρ is defined as the function $\text{grad}_{\text{FR}} \ell_N(\rho) \in L^2(\rho)$ such that

$$\frac{d}{dt} \ell_N(\rho_t) \Big|_{t=0} = g_\rho^{\text{FR}}(\text{grad}_{\text{FR}} \ell_N(\rho), \zeta).$$

To compute it, observe that the right-hand side of the above identity is given by

$$\begin{aligned} \frac{d}{dt} \ell_N(\rho_t) \Big|_{t=0} &= \int \delta \ell_N(\rho) \cdot \partial_t \rho_t \Big|_{t=0} \\ &= \int \delta \ell_N(\rho), \zeta = \int \delta \ell_N(\rho) \left(\alpha - \int \alpha d\rho \right) d\rho \\ &= \int \left(\delta \ell_N(\rho) - \int \delta \ell_N(\rho) d\rho \right) \left(\alpha - \int \alpha d\rho \right) d\rho \\ &= g_\rho^{\text{FR}} \left(\rho \left(\delta \ell_N(\rho) - \int \delta \ell_N(\rho) d\rho \right), \zeta \right). \end{aligned}$$

Therefore

$$g_\rho^{\text{FR}}(\text{grad}_{\text{FR}} \ell_N(\rho), \zeta) = g_\rho^{\text{FR}} \left(\rho \left(\delta \ell_N(\rho) - \int \delta \ell_N(\rho) d\rho \right), \zeta \right)$$

holds for any $\zeta \in \text{Tan}_\rho^{\text{FR}} \mathcal{P}_2(\mathbb{R}^d)$, and as a result

$$\text{grad}_{\text{FR}} \ell_N(\rho) = \rho \left[\delta \ell_N(\rho) - \int \delta \ell_N(\rho) d\rho \right] = \rho [\delta \ell_N(\rho) + 1],$$

where we have used the fact that $\int \delta \ell_N(\rho) d\rho = -1$. Hence, the gradient flow of ℓ_N with respect to the Fisher-Rao metric d_{FR} is given by

$$\partial_t \rho_t = -\text{grad}_{\text{FR}} \ell_N(\rho_t) = -\rho_t [\delta \ell_N(\rho_t) + 1].$$

C.2.2 Other perspectives of Fisher-Rao gradient flow

In this section, we formally illustrate the connection between Fisher-Rao gradient flow (4.1) with proximal gradient descent and mirror descent. For simplicity, we focus on the case when ρ_t is continuous; the case when ρ_t is discrete is similar. We also show the connection between the particle Fisher-Rao gradient descent (2) and the EM algorithm.

Fisher-Rao gradient flow as proximal gradient flow. Consider the proximal gradient update in (4.3). Recall that for $\mu, \nu \in \mathcal{P}_{\text{ac}}(\mathbb{R}^d)$, the Fisher-Rao distance can be expressed as

$$d_{\text{FR}}^2(\mu, \nu) = 4 \int |\sqrt{\mu(x)} - \sqrt{\nu(x)}|^2 dx.$$

Note that for the purpose of defining a gradient flow, the metric only matters up to its second-order local expansion

$$d_{\text{FR}}^2(\mu, \nu) = \int \frac{\Delta^2(x)}{\nu(x)} dx + \text{higher-order terms},$$

where $\Delta = \mu - \nu$. Therefore we obtain an asymptotically (as $\eta \rightarrow 0$) equivalent problem

$$\rho_t^\eta := \arg \min_{\rho \in \mathcal{P}_{\text{ac}}(\mathbb{R}^d)} \left\{ \int_{\mathbb{R}^d} \delta \ell_N(\rho_t) d(\rho - \rho_t) + \frac{1}{2\eta} \int \frac{[\rho(x) - \rho_t(x)]^2}{\rho_t(x)} dx \right\}.$$

The first-order optimality condition is

$$\delta \ell_N(\rho_t)(x) + \frac{1}{\eta} \cdot \frac{\rho(x) - \rho_t(x)}{\rho_t(x)} = c$$

for some constant $c \in \mathbb{R}$. This gives

$$\rho_t^\eta(x) = \rho_t(x) [1 + x\eta - \eta \delta \ell_N(\rho_t)(x)].$$

Since ρ_t^η is a probability density, we have

$$1 = \int_{\mathbb{R}^d} \rho_t^\eta(x) dx = \int_{\mathbb{R}^d} \rho_t(x) [1 + c\eta - \eta \delta \ell_N(\rho_t)(x)] dx = 1 + c\eta - \eta,$$

where we use the fact that $\int_{\mathbb{R}^d} \delta \ell_N(\rho) d\rho = -1$ for any $\rho \in \mathcal{P}_2(\mathbb{R}^d)$. This gives $c = -1$. Therefore

$$\rho_t^\eta(x) = \rho_t(x) [1 - \eta - \eta \delta \ell_N(\rho_t)(x)],$$

and as a result,

$$\partial_t \rho_t = \lim_{\eta \rightarrow 0+} \frac{\rho_t^\eta - \rho_t}{\eta} = -[1 + \delta \ell_N(\rho_t)],$$

which recovers the Fisher-Rao gradient flow (4.1).

Fisher-Rao gradient flow as mirror flow. Recall that the mirror descent update is defined as

$$\rho_t^\eta := \arg \min_{\rho \in \mathcal{P}_{\text{ac}}(\mathbb{R}^d)} \int_{\mathbb{R}^d} \delta \ell_N(\rho_t) d(\rho - \rho_t) + \frac{1}{\eta} \text{KL}(\rho \| \rho_t).$$

The first variation of $f(\cdot) := \text{KL}(\cdot \| \rho_t)$ is given by

$$\delta f(\rho)(x) = \log \left[\frac{\rho(x)}{\rho_t(x)} \right],$$

therefore the first-order optimality condition reads

$$\delta \ell_N(\rho_t)(x) + \frac{1}{\eta} \log \left[\frac{\rho(x)}{\rho_t(x)} \right] = c$$

for some constant $c > 0$. This gives

$$\frac{\rho(x)}{\rho_t(x)} = \exp\{\eta[c - \delta\ell_N(\rho_t)(x)]\} \propto \exp[-\eta\delta\ell_N(\rho_t)(x)].$$

Since $\int_{\mathbb{R}^d} \rho_t^\eta(x) dx = 1$, we know that the closed-form solution is given by

$$\rho_t^\eta(x) = \frac{\rho_t(x) \exp[-\eta\delta\ell_N(\rho_t)(x)]}{\int \rho_t(y) \exp[-\eta\delta\ell_N(\rho_t)(y)] dy}. \quad (\text{C.1})$$

Then as $\eta \rightarrow 0$, we can compute

$$\begin{aligned} \rho_t^\eta(x) &= \frac{\rho_t(x) [1 - \eta\delta\ell_N(\rho_t)(x) + O(\eta^2)]}{\int \rho_t(y) [1 - \eta\delta\ell_N(\rho_t)(y) + O(\eta^2)] dy} = \frac{\rho_t(x) [1 - \eta\delta\ell_N(\rho_t)(x) + O(\eta^2)]}{1 + \eta + O(\eta^2)} \\ &= \rho_t(x) \{1 - \eta[1 + \delta\ell_N(\rho_t)(x)] + O(\eta^2)\}, \end{aligned}$$

where we use the fact that $\int \delta\ell_N(\rho) d\rho = -1$ for any $\rho \in \mathcal{P}_2(\mathbb{R}^d)$. Therefore the continuous-time limit of mirror descent is

$$\partial_t \rho_t = \lim_{\eta \rightarrow 0+} \frac{\rho_t^\eta(x) - \rho_t(x)}{\eta} = -[1 + \delta\ell_N(\rho_t)(x)],$$

which recovers the Fisher-Rao gradient flow (4.1).

Fisher-Rao gradient descent as EM algorithm. Now we consider fitting a m -component Gaussian mixture model with unknown weights $\{\omega^{(j)}\}_{1 \leq j \leq m}$, known location parameters $\{\mu_j\}_{1 \leq j \leq m}$ and isotropic covariance. Given the data $\{X_i\}_{1 \leq i \leq N}$, the MLE is given by

$$\arg \max_{\omega \in \Delta^{m-1}} \ell(\omega) = \frac{1}{N} \sum_{i=1}^N \log \left[\sum_{j=1}^m \omega^{(j)} \phi(X_i - \mu_j) \right].$$

The Expectation-Maximization algorithm for solving the above MLE is given as follows. We first introduce the latent i.i.d. random variables $\{J_i\}_{1 \leq i \leq N}$ distributed $\mathbb{P}(J_i = j) = \omega^{(j)}$ for $1 \leq j \leq m$, then the distribution of the observed samples is $X_i | J_i = j \sim \mathcal{N}(\mu_j, I_d)$. The joint distribution of (X_i, J_i) is given by

$$p_\omega(x, j) = \phi(X_i - \mu_j) \omega^{(j)},$$

and conditional on $X_i = x$, the conditional distribution of J_i is given by

$$p_\omega(j|x) = \frac{p_\omega(x, j)}{\sum_{l=1}^m p_\omega(x, l)} = \frac{\phi(X_i - \mu_j) \omega^{(j)}}{\sum_{l=1}^m \phi(X_i - \mu_l) \omega^{(l)}}.$$

Given the current estimate ω_t , the E-step amounts to computing

$$\begin{aligned} Q(\omega|\omega_t) &= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^m p_{\omega_t}(j|X_i) \log p_\omega(X_i, j) \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^m \frac{\phi(X_i - \mu_j) \omega_t^{(j)}}{\sum_{l=1}^m \phi(X_i - \mu_l) \omega_t^{(l)}} \log [\phi(X_i - \mu_j) \omega^{(j)}]. \end{aligned}$$

The M-step is to update

$$\omega_{t+1} := \arg \max_{\omega \in \Delta^{m-1}} Q(\omega|\omega_t),$$

which is given by

$$\omega_{t+1}^{(j)} = \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - \mu_j) \omega_t^{(j)}}{\sum_{l=1}^m \phi(X_i - \mu_l) \omega_t^{(l)}} \quad \forall 1 \leq j \leq m.$$

This is equivalent to Algorithm 2 with step size $\eta = 1$.

C.3 Wasserstein gradient flow

We introduce the Riemannian structure over $\mathcal{P}_2(\mathbb{R}^d)$ underlying the quadratic Wasserstein distance. We define the tangent space at $\rho \in \mathcal{P}_2(\mathbb{R}^d)$ to be

$$\text{Tan}_\rho^{\text{W}}\mathcal{P}_2(\mathbb{R}^d) = \left\{ \zeta : \zeta = -\text{div}(\rho \nabla u) \text{ for some } u \text{ satisfying } \int \|\nabla u\|_2^2 d\rho < \infty \right\}.$$

We equip this tangent space with the $L^2(\rho)$ metric, namely we define the Riemannian metric tensor $g_\rho^{\text{W}}(\cdot, \cdot) : \text{Tan}_\rho^{\text{W}}\mathcal{P}_2(\mathbb{R}^d) \times \text{Tan}_\rho^{\text{W}}\mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ as

$$g_\rho^{\text{W}}(\zeta_1, \zeta_2) := \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \langle \nabla u_1, \nabla u_2 \rangle \rho(dx)$$

for any $\zeta_1 = -\text{div}(\rho \nabla u_1)$ and $\zeta_2 = -\text{div}(\rho \nabla u_2)$. The metric induced by this Riemannian structure recovers the quadratic Wasserstein distance, namely

$$\begin{aligned} d_{\text{W}}^2(\rho_0, \rho_1) &= \inf \left\{ \int_0^1 \int \|v_t\|_2^2 d\rho_t dt : (\rho_t, v_t)_{t \in [0,1]} \text{ solves } \partial_t \rho_t + \text{div}(\rho_t v_t) = 0 \right\} \\ &= \inf_{\pi \in \Pi(\rho_0, \rho_1)} \int \|x - y\|_2^2 \pi(dx, dy), \end{aligned}$$

where $\Pi(\rho_0, \rho_1)$ is the set of couplings of ρ_0 and ρ_1 . This is known as the Benamou-Brenier formula for the Wasserstein distance. Then we derive the Wasserstein gradient flow with respect to the functional ℓ_N . Interested readers are referred to [Ambrosio et al. \(2008\)](#) for detailed introduction to Wasserstein gradient flow. Let $(\rho_t)_{t \geq 0}$ be a C^1 curve satisfying $\rho_0 = \rho$ with initial velocity

$$\partial_t \rho_t|_{t=0} = \zeta = -\text{div}(\rho \nabla u).$$

Then it should hold that

$$\frac{d}{dt} \ell_N(\rho_t) \Big|_{t=0} = g_\rho^{\text{W}}(\text{grad}_{\text{W}} \ell_N(\rho), \zeta).$$

The left hand side of the above equation equals to

$$\begin{aligned} \frac{d}{dt} \ell_N(\rho_t) \Big|_{t=0} &= \int \delta \ell_N(\rho) \zeta dx = - \int \delta \ell_N(\rho) \text{div}(\rho \nabla u) dx \\ &= - \int \langle \nabla \delta \ell_N(\rho), \nabla u \rangle d\rho \\ &= g_\rho^{\text{W}}(-\text{div}(\nabla \delta \ell_N(\rho) \rho), \zeta). \end{aligned}$$

Therefore

$$g_\rho^{\text{W}}(\text{grad}_{\text{W}} \ell_N(\rho), \zeta) = g_\rho^{\text{W}}(-\text{div}(\nabla \delta \ell_N(\rho) \rho), \zeta)$$

holds for any $\zeta \in \text{Tan}_\rho^{\text{W}}\mathcal{P}_2(\mathbb{R}^d)$, and as a result

$$\text{grad}_{\text{W}} \ell_N(\rho) = -\text{div}(\nabla \delta \ell_N(\rho) \rho).$$

This shows that the gradient flow of ℓ_N with respect to the quadratic Wasserstein distance d_{W} is given by

$$\partial_t \rho_t = -\text{grad}_{\text{W}} \ell_N(\rho_t) = \text{div}(\nabla \delta \ell_N(\rho_t) \rho_t).$$

C.4 Wasserstein-Fisher-Rao gradient flow

We introduce the Riemannian structure over $\mathcal{P}_2(\mathbb{R}^d)$ underlying the Wasserstein-Fisher-Rao metric. Define the tangent space at $\rho \in \mathcal{P}_2(\mathbb{R}^d)$ to be

$$\text{Tan}_\rho^{\text{WFR}}\mathcal{P}_2(\mathbb{R}^d) = \left\{ \zeta : \zeta = -\text{div}(\rho \nabla u) + \rho \left(\alpha - \int \alpha d\rho \right) \text{ for some } u, \alpha : \mathbb{R}^d \rightarrow \mathbb{R} \right\}$$

$$\text{satisfying } \int (\alpha^2 + \|\nabla u\|_2^2) d\rho < \infty \}.$$

We equip this tangent space with the Riemannian metric tensor $g_\rho^{\text{WFR}}(\cdot, \cdot) : \text{Tan}_\rho^{\text{WFR}}\mathcal{P}_2(\mathbb{R}^d) \times \text{Tan}_\rho^{\text{WFR}}\mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ defined as

$$\begin{aligned} g_\rho^{\text{WFR}}(\zeta_1, \zeta_2) &:= \int_{\mathbb{R}^d} \langle \nabla u_1, \nabla u_2 \rangle \rho(dx) + \int_{\mathbb{R}^d} \left[\alpha_1(x) - \int_{\mathbb{R}^d} \alpha_1 d\rho \right] \left[\alpha_2(x) - \int_{\mathbb{R}^d} \alpha_2 d\rho \right] \rho(dx) \\ &= \int_{\mathbb{R}^d} \langle \nabla u_1, \nabla u_2 \rangle \rho(dx) + \int_{\mathbb{R}^d} \alpha_1(x) \alpha_2(x) \rho(dx) - \left(\int_{\mathbb{R}^d} \alpha_1 d\rho \right) \left(\int_{\mathbb{R}^d} \alpha_2 d\rho \right) \end{aligned}$$

for any $\zeta_1 = -\text{div}(\rho \nabla u_1) + \rho(\alpha_1 - \int \alpha_1 d\rho)$ and $\zeta_2 = -\text{div}(\rho \nabla u_2) + \rho(\alpha_2 - \int \alpha_2 d\rho)$. The metric induced by the above Riemannian structure, namely the Wasserstein-Fisher-Rao metric $d_{\text{WFR}}(\cdot, \cdot)$, is defined as

$$\begin{aligned} d_{\text{WFR}}^2(\rho_0, \rho_1) &= \inf \left\{ \int_0^1 \int \left[\|v_t\|^2 + \left(\alpha_t - \int \alpha_t d\rho_t \right)^2 \right] d\rho_t dt : (\rho_t, v_t, \alpha_t)_{0 \leq t \leq 1} \right. \\ &\quad \left. \text{solves } \partial_t \rho_t = -\text{div}(\rho_t v_t) + \rho_t \alpha_t \right\}. \end{aligned}$$

Then we follow [Gallouët and Monsaingeon \(2017\)](#); [Lu et al. \(2019\)](#) to derive the Wasserstein-Fisher-Rao gradient flow with respect to the functional ℓ_N . Let $(\rho_t)_{t \geq 0}$ be a C^1 curve satisfying $\rho_0 = \rho$ with initial velocity

$$\partial_t \rho_t|_{t=0} = \zeta = -\text{div}(\rho \nabla u) + \rho \left(\alpha - \int \alpha d\rho \right).$$

Then it should hold that

$$\frac{d}{dt} \ell_N(\rho_t) \Big|_{t=0} = g_\rho^{\text{WFR}}(\text{grad}_{\text{WFR}} \ell_N(\rho), \zeta).$$

The left hand side of the above equation equals to

$$\begin{aligned} \frac{d}{dt} \ell_N(\rho_t) \Big|_{t=0} &= \int \delta \ell_N(\rho) \zeta dx = \int \delta \ell_N(\rho) \left[-\text{div}(\rho \nabla u) + \rho \left(\alpha - \int \alpha d\rho \right) \right] dx \\ &= - \int \langle \nabla \delta \ell_N(\rho), \nabla u \rangle d\rho + \int \left(\delta \ell_N(\rho) - \int \delta \ell_N(\rho) d\rho \right) \left(\alpha - \int \alpha d\rho \right) d\rho \\ &= g_\rho^{\text{WFR}} \left(-\text{div}(\nabla \delta \ell_N(\rho) \rho) + \rho \left(\delta \ell_N(\rho) - \int \delta \ell_N(\rho) d\rho \right), \zeta \right). \end{aligned}$$

Therefore

$$g_\rho^{\text{WFR}}(\text{grad}_{\text{W}} \ell_N(\rho), \zeta) = g_\rho^{\text{WFR}} \left(-\text{div}(\nabla \delta \ell_N(\rho) \rho) + \rho \left(\delta \ell_N(\rho) - \int \delta \ell_N(\rho) d\rho \right), \zeta \right)$$

holds for any $\zeta \in \text{Tan}_\rho^{\text{WFR}}\mathcal{P}_2(\mathbb{R}^d)$, and as a result

$$\begin{aligned} \text{grad}_{\text{WFR}} \ell_N(\rho) &= -\text{div}(\nabla \delta \ell_N(\rho) \rho) + \rho \left(\delta \ell_N(\rho) - \int \delta \ell_N(\rho) d\rho \right) \\ &= -\text{div}(\nabla \delta \ell_N(\rho) \rho) + \rho [1 + \delta \ell_N(\rho)], \end{aligned}$$

where we have used the fact that $\int \delta \ell_N(\rho) d\rho = -1$. This shows that the gradient flow of ℓ_N with respect to the Wasserstein-Fisher-Rao metric d_{WFR} is given by

$$\partial_t \rho_t = -\text{grad}_{\text{WFR}} \ell_N(\rho_t) = \text{div}(\nabla \delta \ell_N(\rho_t) \rho_t) - \rho_t [1 + \delta \ell_N(\rho_t)].$$

D Proofs for the convergence theory

D.1 Proof of Lemma 1

We invoke a descent lemma Wasserstein gradient descent, whose proof is deferred to Appendix D.2. Such a lemma is standard in convex optimization optimization (see, e.g., Bubeck, 2015, eq. (3.5)). It has appeared for optimization over the Wasserstein space in Salim et al. (2020) under for functionals that are convex along generalized geodesics, an assumption that does not hold for the negative log-likelihood ℓ_N .

Lemma 3 (A descent lemma). *Let Assumption 1 hold. Choose any $\rho \in \mathcal{P}(\mathbb{R}^d)$. Define $c = \inf_{x \in \Omega} (\rho * \phi)(x)$, $G = \sup_{x \in \mathbb{R}^d} \|\nabla \phi(x)\|_2$ and $H = \sup_{x \in \mathbb{R}^d} \|\nabla^2 \phi(x)\|_2$.*

- We have

$$\ell_N(\rho^{W,\eta}) - \ell_N(\rho) \leq -\eta \left[1 - \frac{\eta}{2c} \left(H + \frac{G^2}{c} \right) \right] \mathbb{E}_{Y \sim \rho} \|\nabla \delta \ell_N(\rho)(Y)\|_2^2.$$

In addition, we have $\sup_{x \in \mathbb{R}^d} \|\nabla^2 \delta \ell_N(\rho)(x)\|_2 \leq H/c$.

- If $0 \leq \eta < c/H$ and $\text{supp}(\rho) = \mathbb{R}^d$, then $\text{supp}(\rho^{W,\eta}) = \mathbb{R}^d$.

We now come back to Lemma 1. Let $R = \inf\{r \geq 0 : \bar{\phi}(r) \leq e^{-\ell_N(\rho_0)}/2\}$. Lemma 2 implies that for any $\mu \in \mathcal{P}(\mathbb{R}^d)$ with $\ell_N(\mu) \leq \ell_N(\rho_0)$, we have

$$\inf_{x \in \Omega} (\mu * \phi)(x) \geq e^{-\ell_N(\rho_0)} \underline{\phi}(R + \text{diam}(\Omega))/[2\bar{\phi}(0)] = c_0.$$

In particular, $\inf_{x \in \Omega} (\rho_0 * \phi)(x) \geq c_0$. When $\eta \leq c_0/(H + G^2/c_0)$, Lemma 3 and the definition $\rho_1 = (\text{id} - \eta \nabla \delta \ell_N(\rho_0))_{\#} \rho_0$ together yield

$$\ell_N(\rho_1) - \ell_N(\rho_0) \leq -\frac{\eta}{2} \mathbb{E}_{Y \sim \rho_0} \|\nabla \delta \ell_N(\rho_0)(Y)\|_2^2 \leq -\frac{1}{2\eta} W_2^2(\rho_1, \rho_0) \leq 0.$$

Also, if $\text{supp}(\rho_0) = \mathbb{R}^d$, then $\text{supp}(\rho_1) = \mathbb{R}^d$. From $\ell_N(\rho_1) \leq \ell_N(\rho_0)$ we obtain that $\inf_{x \in \Omega} (\rho_1 * \phi)(x) \geq c_0$. Then, the proof is completed by induction.

D.2 Proof of Lemma 3

We prove the two results in Lemma 3 in sequence.

Part 1. Let $\nu = N^{-1} \sum_{i=1}^N \delta_{X_i}$ be the empirical data distribution, and let $h(x) = -\log x$ for $x > 0$. Then we can write

$$\ell_N(\rho) = -\frac{1}{N} \sum_{i=1}^N \log [\rho * \phi(X_i)] = \mathbb{E}_{X \sim \nu} [h((\rho * \phi)(X))] = \mathbb{E}_{X \sim \nu} [h(\mathbb{E}_{Y \sim \rho} [\phi(X - Y)])]$$

as well as

$$\ell_N(\rho^{W,\eta}) = \mathbb{E}_{X \sim \nu} [h(\mathbb{E}_{Y \sim \rho^{W,\eta}} [\phi(X - Y)])] = \mathbb{E}_{X \sim \nu} [h(\mathbb{E}_{Y \sim \rho} [\phi(X - Y + \eta \nabla \delta \ell_N(\rho)(Y))])].$$

Define $a(x) = \mathbb{E}_{Y \sim \rho} [\phi(x - Y + \eta \nabla \delta \ell_N(\rho)(Y))]$ and $b(x) = \mathbb{E}_{Y \sim \rho} [\phi(x - Y)] = (\rho * \phi)(x)$ for any $x \in \mathbb{R}^d$. Then we can write

$$\ell_N(\rho) = \mathbb{E}_{X \sim \nu} [h(b(X))], \quad \ell_N(\rho^{W,\eta}) = \mathbb{E}_{X \sim \nu} [h(a(X))].$$

Note that $h'(x) = -x^{-1} < 0$, $h''(x) = x^{-2} > 0$ and $h'''(x) = -2x^{-3} < 0$. For any $a, b > 0$, by Taylor's theorem

$$h(a) \leq h(b) + h'(b)(a - b) + \frac{h''(b)}{2}(a - b)^2 = h(b) - \frac{a - b}{b} + \frac{(a - b)^2}{2b^2}.$$

Taking the above two equations collectively gives

$$\begin{aligned}\ell_N(\rho^{W,\eta}) - \ell_N(\rho) &= \mathbb{E}_{X \sim \nu} [h(a(X)) - h(b(X))] \\ &\leq \underbrace{-\mathbb{E}_{X \sim \nu} \left[\frac{a(X) - b(X)}{b(X)} \right]}_{=:\alpha_1} + \underbrace{\frac{1}{2} \mathbb{E}_{X \sim \nu} \left[\frac{[a(X) - b(X)]^2}{b(X)^2} \right]}_{=:\alpha_2}.\end{aligned}\quad (\text{D.1})$$

Now derive upper bounds for α_1 and α_2 respectively.

- To control α_1 , let $G = \sup_{x \in \mathbb{R}^d} \|\nabla \phi(x)\|_2$ and observe that

$$|\phi(x - y + \eta \nabla \delta \ell_N(\rho)(y)) - \phi(x - y)| \leq G \eta \|\nabla \delta \ell_N(\rho)(y)\|_2$$

for all $x, y \in \mathbb{R}^d$, therefore

$$|a(x) - b(x)| \leq G \eta \mathbb{E}_{Y \sim \rho} [\|\nabla \delta \ell_N(\rho)(Y)\|_2] \quad (\text{D.2})$$

holds for all $x \in \mathbb{R}^d$. Then we have

$$\alpha_1 \leq \frac{G^2 \eta^2 \mathbb{E}_{Y \sim \rho} [\|\nabla \delta \ell_N(\rho)(Y)\|_2]}{c^2} \stackrel{\text{(ii)}}{\leq} \frac{G^2 \eta^2}{c^2} \mathbb{E}_{Y \sim \rho} [\|\nabla \delta \ell_N(\rho)(Y)\|_2^2], \quad (\text{D.3})$$

where (i) follows from (D.2) and the fact

$$c = \inf_{x \in \Omega} (\rho * \phi)(x) = \inf_{x \in \Omega} b(x) \quad (\text{D.4})$$

and (ii) follows from Jensen's inequality.

- Regarding α_2 , let $H = \sup_{x \in \mathbb{R}^d} \|\nabla^2 \phi(x)\|_2$ and we have

$$\phi(x - y + \eta \nabla \delta \ell_N(\rho)(y)) - \phi(x - y) \geq \langle \nabla \phi(x - y), \eta \nabla \delta \ell_N(\rho)(y) \rangle - \frac{H}{2} \eta^2 \|\nabla \delta \ell_N(\rho)(y)\|_2^2,$$

for any $x, y \in \mathbb{R}^d$, and therefore

$$a(x) - b(x) \geq \mathbb{E}_{Y \sim \rho} [\langle \nabla \phi(x - Y), \eta \nabla \delta \ell_N(\rho)(Y) \rangle] - \frac{H}{2} \eta^2 \mathbb{E}_{Y \sim \rho} [\|\nabla \delta \ell_N(\rho)(Y)\|_2^2] \quad (\text{D.5})$$

for any $x \in \mathbb{R}^d$. Since $b(x) > 0$, we have

$$\begin{aligned}\alpha_2 &\stackrel{\text{(i)}}{\leq} -\mathbb{E}_{X \sim \nu} \left[\frac{\mathbb{E}_{Y \sim \rho} [\langle \nabla \phi(X - Y), \eta \nabla \delta \ell_N(\rho)(Y) \rangle] - \frac{H}{2} \eta^2 \mathbb{E}_{Y \sim \rho} [\|\nabla \delta \ell_N(\rho)(Y)\|_2^2]}{b(X)} \right] \\ &= -\eta \mathbb{E}_{Y \sim \rho} \left[\left\langle \mathbb{E}_{X \sim \nu} \left(\frac{\nabla \phi(X - Y)}{b(X)} \right), \nabla \delta \ell_N(\rho)(Y) \right\rangle \right] + \frac{H \eta^2}{2} \mathbb{E}_{X \sim \nu} \left[\frac{1}{b(X)} \right] \mathbb{E}_{Y \sim \rho} [\|\nabla \delta \ell_N(\rho)(Y)\|_2^2] \\ &\stackrel{\text{(ii)}}{=} -\eta \mathbb{E}_{Y \sim \rho} [\|\nabla \delta \ell_N(\rho)(Y)\|_2^2] + \frac{H \eta^2}{2} \mathbb{E}_{X \sim \nu} \left[\frac{1}{b(X)} \right] \mathbb{E}_{Y \sim \rho} [\|\nabla \delta \ell_N(\rho)(Y)\|_2^2] \\ &\stackrel{\text{(iii)}}{\leq} \left(-\eta + \frac{H \eta^2}{2c} \right) \mathbb{E}_{Y \sim \rho} [\|\nabla \delta \ell_N(\rho)(Y)\|_2^2].\end{aligned}\quad (\text{D.6})$$

Here (i) utilizes (D.5); (ii) holds since

$$\delta \ell_N(\rho)(y) = -\mathbb{E}_{X \sim \nu} \left[\frac{\phi(X - y)}{b(X)} \right], \quad \nabla \delta \ell_N(\rho)(y) = \mathbb{E}_{X \sim \nu} \left[\frac{\nabla \phi(X - y)}{b(X)} \right], \quad (\text{D.7})$$

and therefore

$$\mathbb{E}_{Y \sim \rho} \left[\left\langle \mathbb{E}_{X \sim \nu} \left(\frac{\nabla \phi(X - Y)}{b(X)} \right), \nabla \delta \ell_N(\rho)(Y) \right\rangle \right] = \mathbb{E}_{Y \sim \rho} [\|\nabla \delta \ell_N(\rho)(Y)\|_2^2];$$

and (iii) follows from (D.4).

Taking (D.1), (D.3) and (D.6) collectively gives

$$\begin{aligned}\ell_N(\rho^{\text{W},\eta}) - \ell_N(\rho) &\leq \alpha_1 + \alpha_2 \leq \left(-\eta + \frac{H\eta^2}{2c} + \frac{G^2\eta^2}{2c^2}\right) \mathbb{E}_{Y \sim \rho} \left[\|\nabla \delta \ell_N(\rho)(Y)\|_2^2\right] \\ &= -\eta \left[1 - \frac{\eta}{2c} \left(H + \frac{G^2}{c}\right)\right] \mathbb{E}_{Y \sim \rho} \left[\|\nabla \delta \ell_N(\rho)(Y)\|_2^2\right].\end{aligned}$$

Finally, we learn from (D.7) and (D.4) that

$$\|\nabla^2 \delta \ell_N(\rho)(x)\|_2 = \left\| \mathbb{E}_{X \sim \nu} \left[\frac{\nabla^2 \phi(X - x)}{b(X)} \right] \right\|_2 \leq \frac{\sup_{x \in \mathbb{R}^d} \|\nabla^2 \phi(x)\|_2}{\inf_{x \in \mathbb{R}^d} b(x)} = \frac{H}{c}.$$

Part 2. When $\eta < c/H$, the mapping $x \mapsto \eta \nabla \delta \ell_N(\rho)(x)$ is a contraction. Let $\varphi(x) = x - \eta \nabla \delta \ell_N(\rho)(x)$. By Lemma 3, $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a bijection and φ^{-1} is Lipschitz. The second-order differentiability of $\delta \ell_N(\rho)$ implies the differentiability of φ and thus φ^{-1} . If $\text{supp}(\rho) = \mathbb{R}^d$, then $\text{supp}(\rho^{\text{W},\eta}) = \text{supp}(\varphi_{\#}\rho) = \mathbb{R}^d$.

D.3 Proof of Lemma 2

Let $\nu = N^{-1} \sum_{i=1}^N \delta_{X_i}$ be the empirical data distribution. For any $\mu \ll \rho$, we have

$$\begin{aligned}\ell_N(\mu) &= -\mathbb{E}_{X \sim \nu} [\log((\mu * \phi)(X))] = -\mathbb{E}_{X \sim \nu} [\log(\mathbb{E}_{Y \sim \mu} [\phi(X - Y)])] \\ &= -\mathbb{E}_{X \sim \nu} \left[\log \left(\mathbb{E}_{Y \sim \rho} \left[\frac{d\mu}{d\rho}(Y) \cdot \phi(X - Y) \right] \right) \right] \\ &= -\mathbb{E}_{X \sim \nu} \left[\log \left(\mathbb{E}_{Y \sim \rho} \left[\frac{d\mu}{d\rho}(Y) \cdot (\rho * \phi)(X) \cdot \frac{\phi(X - Y)}{(\rho * \phi)(X)} \right] \right) \right].\end{aligned}$$

For any $x \in \mathbb{R}^d$, we can define a new probability measure $\rho_x^{\text{FR},1} \in \mathcal{P}(\mathbb{R}^d)$ through

$$\frac{d\rho_x^{\text{FR},1}}{d\rho}(\cdot) = \frac{\phi(x - \cdot)}{(\rho * \phi)(x)}, \quad (\text{D.8})$$

and we can check that $\rho_x^{\text{FR},1}$ is indeed a probability measure since

$$\int_{\mathbb{R}^d} d\rho_x^{\text{FR},1} = \int_{\mathbb{R}^d} \frac{\phi(x - y)}{(\rho * \phi)(x)} \rho(dy) = 1.$$

Then we have, by the convexity of $t \mapsto -\log t$ and Jensen's inequality,

$$\begin{aligned}\ell_N(\mu) &\stackrel{(i)}{=} -\mathbb{E}_{X \sim \nu} \left[\log \left(\mathbb{E}_{Y \sim \rho_X^{\text{FR},1}} \left[\frac{d\mu}{d\rho}(Y) \cdot (\rho * \phi)(X) \right] \right) \right] \\ &\stackrel{(ii)}{\leq} -\mathbb{E}_{X \sim \nu} \left[\mathbb{E}_{Y \sim \rho_X^{\text{FR},1}} \left[\log \left(\frac{d\mu}{d\rho}(Y) \cdot (\rho * \phi)(X) \right) \right] \right] \\ &\stackrel{(iii)}{=} -\mathbb{E}_{X \sim \nu} \left[\mathbb{E}_{Y \sim \rho} \left[\log \left(\frac{d\mu}{d\rho}(Y) \cdot (\rho * \phi)(X) \right) \cdot \frac{\phi(X - Y)}{(\rho * \phi)(X)} \right] \right] \\ &= -\underbrace{\mathbb{E}_{X \sim \nu, Y \sim \rho} \left[\log \left(\frac{d\mu}{d\rho}(Y) \right) \cdot \frac{\phi(X - Y)}{(\rho * \phi)(X)} \right]}_{=:\beta_1} - \underbrace{\mathbb{E}_{X \sim \nu, Y \sim \rho} \left[\log((\rho * \phi)(X)) \cdot \frac{\phi(X - Y)}{(\rho * \phi)(X)} \right]}_{=:\beta_2}. \quad (\text{D.9})\end{aligned}$$

Here (i) and (iii) utilizes (D.8), and (ii) follows from Jensen's inequality and the convexity of $t \mapsto -\log t$. Then we study the two terms β_1 and β_2 respectively. Regarding β_1 , we have

$$\begin{aligned}\beta_1 &= \mathbb{E}_{Y \sim \rho} \left[\log \left(\frac{d\mu}{d\rho}(Y) \right) \mathbb{E}_{X \sim \nu} \left[\frac{\phi(X - Y)}{(\rho * \phi)(X)} \right] \right] \\ &= \mathbb{E}_{Y \sim \rho} \left[\log \left(\frac{d\mu}{d\rho}(Y) \right) \cdot \frac{d\rho_x^{\text{FR},1}}{d\rho}(Y) \right] = \mathbb{E}_{Y \sim \rho^{\text{FR},1}} \left[\log \left(\frac{d\mu}{d\rho}(Y) \right) \right]\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_{Y \sim \rho^{\text{FR},1}} \left[\log \left(\frac{d\mu}{d\rho^{\text{FR},1}}(Y) \right) \right] + \mathbb{E}_{Y \sim \rho^{\text{FR},1}} \left[\log \left(\frac{d\rho^{\text{FR},1}}{d\rho}(Y) \right) \right] \\
&= -\text{KL}(\rho^{\text{FR},1} \parallel \mu) + \text{KL}(\rho^{\text{FR},1} \parallel \rho).
\end{aligned} \tag{D.10}$$

Regarding β_2 , we have

$$\beta_2 = \mathbb{E}_{X \sim \nu} \left[\log((\rho * \phi)(X)) \mathbb{E}_{Y \sim \rho} \left[\frac{\phi(X - Y)}{(\rho * \phi)(X)} \right] \right] = \mathbb{E}_{X \sim \nu} [\log((\rho * \phi)(X))] = -\ell_N(\rho). \tag{D.11}$$

Taking (D.9), (D.10) and (D.11) collectively yields

$$\ell_N(\mu) \leq \ell_N(\rho) + \text{KL}(\rho^{\text{FR},1} \parallel \mu) - \text{KL}(\rho^{\text{FR},1} \parallel \rho), \quad \forall \mu \ll \rho.$$

By taking $\mu = \rho^{\text{FR},1}$, we get

$$\ell_N(\rho^{\text{FR},1}) \leq \ell_N(\rho) - \text{KL}(\rho^{\text{FR},1} \parallel \rho). \tag{D.12}$$

Recall that for any $\gamma \in (0, 1)$ we have $\rho^{\text{FR},\gamma} = (1 - \gamma)\rho + \gamma\rho^{\text{FR},1}$. Therefore

$$\ell_N(\rho^{\text{FR},\gamma}) \stackrel{(i)}{\leq} (1 - \gamma)\ell_N(\rho) + \gamma\ell_N(\rho^{\text{FR},1}) \stackrel{(ii)}{\leq} \ell_N(\rho) - \gamma\text{KL}(\rho^{\text{FR},1} \parallel \rho), \tag{D.13}$$

where (i) holds since $\ell_N(\rho)$ is ℓ_2 -convex in ρ , and (ii) follows from (D.12). By the ℓ_2 -convexity of $\text{KL}(\cdot \parallel \rho)$, we have

$$\text{KL}(\rho^{\text{FR},\gamma} \parallel \rho) \leq (1 - \gamma)\text{KL}(\rho \parallel \rho) + \gamma\text{KL}(\rho^{\text{FR},1} \parallel \rho) = \gamma\text{KL}(\rho^{\text{FR},1} \parallel \rho). \tag{D.14}$$

Combine (D.13) and (D.14) to achieve

$$\ell_N(\rho^{\text{FR},\gamma}) \leq \ell_N(\rho) - \text{KL}(\rho^{\text{FR},\gamma} \parallel \rho).$$

E Well-posedness of particle gradient flows

E.1 Fisher-Rao gradient flow (Proof of Theorem 5)

In this section, we will show that for the ODE system (4.5)

$$\dot{\omega}_t^{(j)} = -\omega_t^{(j)} \left[1 - \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - \mu^{(j)})}{\sum_{l=1}^m \omega_t^{(l)} \phi(X_i - \mu^{(l)})} \right], \quad \forall t \geq 0, j \in [m]$$

with initial value $\omega_0 \in \Delta^{m-1}$, the solution exists, is unique, and $(\rho_t)_{t \geq 0}$ where $\rho_t := \sum_{l=1}^m \omega_t^{(l)} \delta_{\mu^{(l)}}$ is a Fisher-Rao gradient flow in the sense of (4.1).

First of all, we will use Picard-Lindelöf theorem to prove existence and uniqueness of the solution. Define a function $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$ as

$$f(y) = [f_j(y)]_{1 \leq j \leq m}, \quad f_j(y) = -y_j \left[1 - \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - \mu^{(j)})}{\sum_{l=1}^m y_l \phi(X_i - \mu^{(l)})} \right].$$

Then we can rewrite the ODE system as $\dot{\omega}_t = f(\omega_t)$. For any $y \in \mathbb{R}^m$ such that $\|y - \omega_0\|_2 \leq \varepsilon$ for some $\delta > 0$ to be specified later, we know that

$$\begin{aligned}
\sum_{l=1}^m y_l \phi(X_i - \mu^{(l)}) &\geq \sum_{l=1}^m \omega_0^{(l)} \phi(X_i - \mu^{(l)}) - \sum_{l=1}^m (\omega_0^{(l)} - y_l) \phi(X_i - \mu^{(l)}) \\
&\stackrel{(i)}{\geq} \min_{i \in [N], l \in [m]} \phi(X_i - \mu^{(l)}) - \|y - \omega_0\|_2 \sum_{l=1}^m \phi^2(X_i - \mu^{(l)}) \\
&\stackrel{(ii)}{\geq} \min_{i \in [N], l \in [m]} \phi(X_i - \mu^{(l)}) - \frac{\varepsilon m}{(2\pi)^d}
\end{aligned}$$

for any $i \in [N]$, where (i) follows from $\omega_0^{(l)} \in \Delta^{m-1}$ and Cauchy-Schwarz inequality, while (ii) holds since $\|\phi\|_\infty \leq (2\pi)^{-d/2}$. Therefore by taking

$$\varepsilon := \frac{(2\pi)^d}{2m} \min_{i \in [N], l \in [m]} \phi(X_i - \mu^{(l)}),$$

we know that for any $\|y - \omega_0\|_2 \leq \varepsilon$ it holds that

$$\sum_{l=1}^m y_l \phi(X_i - \mu^{(l)}) \geq \frac{1}{2} \min_{i \in [N], l \in [m]} \phi(X_i - \mu^{(l)}) \triangleq \delta. \quad (\text{E.1})$$

Therefore we can check that

$$\begin{aligned} |[\nabla f_j(y)]_j| &= \left| - \left[1 - \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - \mu^{(j)})}{\sum_{l=1}^m y_l \phi(X_i - \mu^{(l)})} \right] - \left[\frac{1}{N} \sum_{i=1}^N \frac{y_j \phi^2(X_i - \mu^{(j)})}{[\sum_{l=1}^m y_l \phi(X_i - \mu^{(l)})]^2} \right] \right| \\ &\leq 1 + \frac{\|\phi\|_\infty}{\delta} + \frac{(1+\varepsilon)\|\phi\|_\infty^2}{\delta^2} = 1 + \frac{1}{(2\pi)^{d/2} \delta} + \frac{1+\varepsilon}{(2\pi)^d \delta^2}, \end{aligned}$$

and for $l \neq j$

$$|[\nabla f_j(y)]_l| = \left| -y_j \left[\frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - \mu^{(j)}) \phi(X_i - \mu^{(l)})}{[\sum_{l=1}^m y_l \phi(X_i - \mu^{(l)})]^2} \right] \right| \leq \frac{(1+\varepsilon)\|\phi\|_\infty^2}{\delta^2} = \frac{1+\varepsilon}{(2\pi)^d \delta^2}.$$

As a result, we know that for any

$$\max_{y: \|y - \omega_0\|_2 \leq \varepsilon} \|\nabla f_j(y)\|_2 \leq \sqrt{m} \left(1 + \frac{1}{(2\pi)^{d/2} \delta} + \frac{1+\varepsilon}{(2\pi)^d \delta^2} \right) \triangleq C_{\text{lip}}, \quad (\text{E.2})$$

and hence $f(y)$ is C_{lip} -Lipschitz continuous in $\{y : \|y - \omega_0\|_2 \leq \varepsilon\}$ where $C_{\text{lip}} := \sqrt{m} C_{\text{lip}}$. In addition, it is easy to show that

$$\max_{y: \|y - \omega_0\|_2 \leq \varepsilon} \|f_j(y)\|_2 \leq -y_j + \frac{1}{N} \sum_{i=1}^N \frac{y_j \phi(X_i - \mu^{(j)})}{\sum_{l=1}^m y_l \phi(X_i - \mu^{(l)})} \leq 1 + \varepsilon \triangleq M.$$

By Picard-Lindelöf theorem, there exists $t_0 > 0$ such that the ODE has a unique solution on the time interval $[0, t_0]$. We first check that $\omega_t^{(j)} > 0$ for any $j \in [m]$ and $t \in [0, t_0]$. If this is not true, define

$$t_\star := \min \left\{ t \in [0, t_0] : \exists j \in [m] \text{ s.t. } \omega_t^{(j)} \leq 0 \right\}$$

and suppose $\omega_{t_\star}^{(j_\star)} \leq 0$ for $j_\star \in [m]$. Then we know that $\omega_t^{(j)} > 0$ for any $j \in [m]$ and $0 \leq t \leq t_\star$, and hence $\dot{\omega}_t^{(j)} \geq -\omega_t^{(j)}$ for all $t \in [0, t_\star]$. Then we can use Grönwall's lemma to achieve $\omega_t^{(j)} \geq \omega_0^{(j)} e^{-t}$ for all $t \in [0, t_\star]$, and as a result $\omega_{t_\star}^{(j_\star)} > 0$, which is a contradiction. In addition, we can also check that

$$\frac{d}{dt} \sum_{j=1}^m \omega_t^{(j)} = \sum_{j=1}^m \dot{\omega}_t^{(j)} = - \sum_{j=1}^m \omega_t^{(j)} \left[1 - \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - \mu^{(j)})}{\sum_{l=1}^m \omega_t^{(l)} \phi(X_i - \mu^{(l)})} \right] = 0$$

for all $t \in [0, t_0]$. As a result $\omega_t \in \Delta^{m-1}$ for any $0 \leq t \leq t_0$. By repeating the same procedure as above (notice that the above proof only depends on $\omega_0 \in \Delta^{m-1}$, and t_0 only depends on universal constants C_{lip} and M and does not depend on ω_0), we can show that the ODE has a unique solution on $[t_0, 2t_0]$, $[2t_0, 3t_0]$, and so on. This shows the existence and uniqueness of the solution to the ODE system (4.5).

Next, we show that $(\rho_t)_{t \geq 0}$ defined as $\rho_t := \sum_{l=1}^m \omega_t^{(l)} \delta_{\mu^{(l)}}$ solves (4.1). Note that ρ_t is a probability measure since we have shown that $\omega_t \in \Delta^{m-1}$ for any $t \geq 0$. For any test function $\varphi(x) \in C_c^\infty$, we have

$$\frac{d}{dt} \int_{\mathbb{R}^d} \varphi(x) \rho_t(dx) = \frac{d}{dt} \left[\sum_{j=1}^m \omega_t^{(j)} \varphi(\mu_t^{(j)}) \right] = \sum_{j=1}^m \dot{\omega}_t^{(j)} \varphi(\mu_t^{(j)})$$

$$\begin{aligned}
&= - \sum_{j=1}^m \left[1 + \delta \ell_N(\rho_t) \left(\mu_t^{(j)} \right) \right] \omega_t^{(j)} \varphi \left(\mu_t^{(j)} \right) \\
&= - \int_{\mathbb{R}^d} [1 + \delta \ell_N(\rho_t)(x)] \varphi(x) \rho_t(dx).
\end{aligned}$$

This proves that

$$\partial_t \rho_t = -[\delta \ell(\rho_t) + 1] \rho_t$$

holds in the sense of distributions.

E.2 Wasserstein gradient flow (Proof of Theorem 6)

In this section we will show that the ODE system (4.9)

$$\dot{\mu}_t^{(j)} = \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - \mu_t^{(j)})}{m^{-1} \sum_{l=1}^m \phi(X_i - \mu_t^{(l)})} (X_i - \mu_t^{(j)}), \quad \forall t \geq 0, j \in [m]$$

has unique solution, and $(\rho_t)_{t \geq 0}$ where $\rho_t := m^{-1} \sum_{l=1}^m \delta_{\mu_t^{(l)}}$ is a Wasserstein gradient flow in the sense of (4.8).

We will use Picard-Lindelöf theorem to prove existence and uniqueness of the solution. Define a sufficiently large constant

$$R := \max_{1 \leq i \leq N} \|X_i\|_2.$$

For each $j \in [m]$, define a function $f^{(j)} : \mathbb{R}^{md} \rightarrow \mathbb{R}^d$ as

$$f^{(j)}(z) = \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - z_j)}{m^{-1} \sum_{l=1}^m \phi(X_i - z_l)} (X_i - z_j),$$

where $z = [z_j]_{1 \leq j \leq m} \in \mathbb{R}^{md}$ and $z_1, \dots, z_m \in \mathbb{R}^d$, and let $f(z) = [f^{(j)}(z)]_{1 \leq j \leq m}$. Then we can write the ODE system as $\dot{\mu}_t = f(\mu_t)$ where $\mu_t = [\mu_t^{(j)}]_{1 \leq j \leq m}$. Denote by $f^{(j)}(z) = [f_k^{(j)}(z)]_{1 \leq k \leq d}$. For any $z \in \mathbb{R}^{md}$ satisfying $\max_{j \in [m]} \|z_j\|_2 \leq 2R$, we have

$$\min_{1 \leq i \leq N} \frac{1}{m} \sum_{l=1}^m \phi(X_i - z_l) \geq \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{9}{2}R^2\right) \triangleq \delta. \quad (\text{E.3})$$

Then we can compute for $l \neq j$

$$\begin{aligned}
\left\| \nabla_{z_l} f_k^{(j)}(z) \right\|_2 &= \left\| \nabla_{z_l} \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - z_j)}{m^{-1} \sum_{l=1}^m \phi(X_i - z_l)} e_k^\top (X_i - z_j) \right\|_2 \\
&= \left\| -\frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - z_j) m^{-1} \phi(X_i - z_l)}{[m^{-1} \sum_{l=1}^m \phi(X_i - z_l)]^2} e_k^\top (X_i - z_j) (X_i - z_l) \right\|_2 \\
&\stackrel{(i)}{\leq} \frac{m^{-1} \|\phi\|_\infty^2}{\delta^2} \frac{1}{N} \sum_{i=1}^N |e_k^\top (X_i - z_j)| \|X_i - z_l\|_2 \stackrel{(ii)}{\leq} \frac{9m^{-1}}{\delta^2 (2\pi)^d} R^2,
\end{aligned}$$

where (i) utilizes (E.3) and (ii) follows from $\max_{i \in [N]} \|X_i\|_2 \leq R$ and $\max_{j \in [m]} \|z_j\|_2 \leq 2R$. Similarly we have

$$\begin{aligned}
\left\| \nabla_{z_j} f_k^{(j)}(z) \right\|_2 &= \left\| \nabla_{z_j} \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - z_j)}{m^{-1} \sum_{l=1}^m \phi(X_i - z_l)} e_k^\top (X_i - z_j) \right\|_2 \\
&\leq \left\| -\frac{1}{N} \sum_{i=1}^N \frac{\phi^2(X_i - z_j) m^{-1}}{[m^{-1} \sum_{l=1}^m \phi(X_i - z_l)]^2} e_k^\top (X_i - z_j) (X_i - z_j) \right\|_2
\end{aligned}$$

$$\begin{aligned}
& + \left\| \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - z_j)}{m^{-1} \sum_{l=1}^m \phi(X_i - z_l)} [e_k^\top (X_i - z_j) (X_i - z_j) - e_k] \right\|_2 \\
& \leq 9 \left(\frac{m^{-1}}{\delta^2 (2\pi)^d} + \frac{1}{\delta (2\pi)^{d/2}} \right) R^2 + \frac{\|\phi\|_\infty}{\delta}.
\end{aligned}$$

As a result, we have

$$\left\| \nabla_z f_k^{(j)}(z) \right\|_2 = \sqrt{\sum_{l=1}^m \left\| \nabla_{z_l} f_k^{(j)}(z) \right\|_2^2} \leq \sqrt{(m+1) \left(\frac{9m^{-1}}{\delta^2 (2\pi)^d} R^2 \right)^2 + 2 \frac{1}{\delta^2 (2\pi)^d}} \triangleq C_{\text{lip}}. \quad (\text{E.4})$$

Therefore $f^{(j)}(z)$ is $\sqrt{d}C_{\text{Lip}}$ -Lipschitz continuous in $\{z : \max_{j \in [m]} \|z_j\|_2 \leq 2R\}$, and hence $f(z)$ is C_{Lip} -Lipschitz continuous where $C_{\text{Lip}} := C_{\text{lip}} \sqrt{md}$. In addition, it is straightforward to show that for any $z \in \mathbb{R}^{md}$ satisfying $\max_{j \in [m]} \|z_j\|_2 \leq 2R$,

$$\left\| f^{(j)}(z) \right\|_2 \leq \frac{\|\phi\|_\infty}{\delta} 3R \triangleq M$$

holds for all $1 \leq j \leq m$.

Recall that $\mu_0^{(1)}, \dots, \mu_0^{(m)}$ are i.i.d. sampled from $\text{Uniform}(\{X_i\}_{1 \leq i \leq N})$, therefore $\max_{j \in [m]} \|\mu_0^{(j)}\|_2 \leq R$. Therefore we have

$$\{z : \|z - \mu_0\|_2 \leq R\} \subseteq \left\{ z : \max_{j \in [m]} \|z_j\|_2 \leq 2R \right\},$$

where $\mu_0 = [\mu_0^{(j)}]_{1 \leq j \leq m}$. Hence $f(z)$ is $\sqrt{md}C_{\text{Lip}}$ -Lipschitz continuous in $\{z : \|z - \mu_0\|_2 \leq R\}$. Then we can use Picard-Lindelöf theorem to show that, there exists $t_0 > 0$ such that the ODE has a unique solution on the time interval $[0, t_0]$. For any $t \in [0, t_0]$ and $j \in [m]$, we can compute

$$\begin{aligned}
\frac{d}{dt} \|\mu_t^{(j)}\|_2^2 &= 2 \langle \mu_t^{(j)}, \dot{\mu}_t^{(j)} \rangle = \frac{2}{N} \sum_{i=1}^N \frac{\phi(X_i - \mu_t^{(j)})}{m^{-1} \sum_{l=1}^m \phi(X_i - \mu_t^{(l)})} \mu_t^{(j)\top} (X_i - \mu_t^{(j)}) \\
&\leq \frac{2}{N} \sum_{i=1}^N \frac{\phi(X_i - \mu_t^{(j)})}{m^{-1} \sum_{l=1}^m \phi(X_i - \mu_t^{(l)})} \left(\|X_i\|_2 \|\mu_t^{(j)}\|_2 - \|\mu_t^{(j)}\|_2^2 \right) \\
&\leq \frac{2}{N} \sum_{i=1}^N \frac{\phi(X_i - \mu_t^{(j)})}{m^{-1} \sum_{l=1}^m \phi(X_i - \mu_t^{(l)})} \left(R - \|\mu_t^{(j)}\|_2 \right) \|\mu_t^{(j)}\|_2,
\end{aligned}$$

where we use Cauchy-Schwarz inequality in the penultimate step. This shows that $\frac{d}{dt} \|\mu_t^{(j)}\|_2^2 < 0$ as long as $\|\mu_{t_0}^{(j)}\|_2 > R$, and as a result $\max_{j \in [m]} \|\mu_{t_0}^{(j)}\|_2 \leq R$. Then we can repeat the same analysis as above (notice that the above proof only requires $\max_{j \in [m]} \|\mu_0^{(j)}\|_2 \leq R$, and t_0 only depends on universal constants C_{Lip} and M) to show that the ODE has a unique solution on $[t_0, 2t_0]$, $[2t_0, 3t_0]$, and so on. This shows the existence and uniqueness of the solution to the ODE system (4.5).

Finally we check that $(\rho_t)_{t \geq 0}$ defined as $\rho_t := \sum_{l=1}^m m^{-1} \delta_{\mu_t^{(l)}}$ solves (4.1). For any test function $\varphi(x) \in C_c^\infty$, we have

$$\begin{aligned}
\frac{d}{dt} \int_{\mathbb{R}^d} \varphi(x) \rho_t(dx) &= \frac{d}{dt} \left[\frac{1}{m} \sum_{j=1}^m \varphi(\mu_t^{(j)}) \right] = \frac{1}{m} \sum_{j=1}^m \langle \nabla \varphi(\mu_t^{(j)}), \dot{\mu}_t^{(j)} \rangle \\
&= \frac{1}{m} \sum_{j=1}^m \langle \nabla \varphi(\mu_t^{(j)}), -\nabla \delta \ell_N(\rho_t)(\mu_t^{(j)}) \rangle \\
&= - \int_{\mathbb{R}^d} \langle \nabla \varphi(x), \nabla \delta \ell_N(\rho_t) \rangle \rho_t(dx).
\end{aligned}$$

This proves that

$$\partial_t \rho_t = \operatorname{div}(\rho_t \nabla \delta \ell(\rho_t))$$

holds in the sense of distributions.

E.3 Wasserstein-Fisher-Rao gradient flow (Proof of Theorem 3)

In this section we will show that the ODE system (3.8)

$$\begin{aligned}\dot{\mu}_t^{(j)} &= \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - \mu_t^{(j)})}{m^{-1} \sum_{l=1}^m \phi(X_i - \mu_t^{(l)})} (X_i - \mu_t^{(j)}), \\ \dot{\omega}_t^{(j)} &= \left[\frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - \mu_t^{(j)})}{\sum_{l=1}^m \omega_t^{(j)} \phi(X_i - \mu_t^{(l)})} - 1 \right] \omega_t^{(j)},\end{aligned}$$

has unique solution, and $(\rho_t)_{t \geq 0}$ where $\rho_t := \sum_{l=1}^m \omega_t^{(l)} \delta_{\mu_t^{(l)}}$ is a Wasserstein-Fisher-Rao gradient flow in the sense of (3.6). We will integrate the proof techniques used in the previous two sections.

We will again use Picard-Lindelöf theorem to prove existence and uniqueness of the solution. For each $j \in [m]$, define two functions $f^{(j)} : \Delta^{m-1} \times \mathbb{R}^{md} \rightarrow \mathbb{R}^d$ and $g^{(j)} : \Delta^{m-1} \times \mathbb{R}^{md} \rightarrow \mathbb{R}$ as

$$\begin{aligned}f^{(j)}(y, z) &= \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - z_j)}{\sum_{l=1}^m y_l \phi(X_i - z_l)} (X_i - z_j), \\ g^{(j)}(y, z) &= - \left[1 - \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - z_j)}{\sum_{l=1}^m y_l \phi(X_i - z_l)} \right] y_j,\end{aligned}$$

where $y = [y_j]_{1 \leq j \leq m} \in \Delta^{m-1}$, $z = [z_j]_{1 \leq j \leq m} \in \mathbb{R}^{md}$ with $z_1, \dots, z_m \in \mathbb{R}^d$. Let

$$f(y, z) := \begin{bmatrix} f^{(1)}(y, z) \\ \vdots \\ f^{(m)}(y, z) \end{bmatrix}, \quad g(y, z) := \begin{bmatrix} g^{(1)}(y, z) \\ \vdots \\ g^{(m)}(y, z) \end{bmatrix}, \quad \text{and} \quad h(y, z) = \begin{bmatrix} f(y, z) \\ g(y, z) \end{bmatrix}.$$

Then we can write the ODE system as

$$\begin{bmatrix} \dot{\mu}_t \\ \dot{\omega}_t \end{bmatrix} = h \left(\begin{bmatrix} \mu_t \\ \omega_t \end{bmatrix} \right),$$

where $\mu_t = [\mu_t^{(j)}]_{1 \leq j \leq m}$ and $\omega_t = [\omega_t^{(j)}]_{1 \leq j \leq m}$. Denote by $f^{(j)}(z) = [f_k^{(j)}(z)]_{1 \leq k \leq d}$.

For any $z \in \mathbb{R}^{md}$ satisfying $\max_{j \in [m]} \|z_j\|_2 \leq 2R$ and any $y \in \mathbb{R}^m$ satisfying $\|y - \omega_0\|_2 \leq \varepsilon$ where

$$R := \max_{1 \leq i \leq N} \|X_i\|_2, \quad \varepsilon := \frac{(2\pi)^{d/2}}{2m} \exp\left(-\frac{9}{2}R^2\right),$$

we have for any $i \in [N]$

$$\begin{aligned}\sum_{l=1}^m y_l \phi(X_i - z_l) &\geq \sum_{l=1}^m \omega_0^{(l)} \phi(X_i - z_l) - \sum_{l=1}^m (\omega_0^{(l)} - y_l) \phi(X_i - z_l) \\ &\stackrel{(i)}{\geq} \min_{i \in [N], l \in [m]} \phi(X_i - z_l) - \|y - \omega_0\|_2 \sum_{l=1}^m \phi^2(X_i - \mu^{(l)}) \\ &\stackrel{(ii)}{\geq} \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{9}{2}R^2\right) - \frac{\varepsilon m}{(2\pi)^d} = \frac{1}{2(2\pi)^{d/2}} \exp\left(-\frac{9}{2}R^2\right) \triangleq \delta.\end{aligned} \tag{E.5}$$

Here (i) follows from $\omega_0 \in \Delta^{m-1}$ and the Cauchy-Schwarz inequality, while (ii) and (iii) holds since $\|X_i - z_l\|_2 \leq 3R$ for any $i \in [N]$ and $l \in [m]$. For any $j \in [m]$, denote by $f^{(j)}(z) = [f_k^{(j)}(z)]_{1 \leq k \leq d}$. Similar to the proof of (E.4), we can use (E.5) to show that

$$\left\| \nabla_z f_k^{(j)}(y, z) \right\|_2 \leq \sqrt{(m+1) \left(\frac{9m^{-1}}{\delta^2 (2\pi)^d} R^2 \right)^2 + 2 \frac{1}{\delta^2 (2\pi)^d}}.$$

We also have

$$\begin{aligned} \nabla_y f_k^{(j)}(y, z) &= \nabla_y \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - z_j)}{\sum_{l=1}^m y_l \phi(X_i - z_l)} e_k^\top (X_i - z_j) \\ &= -\frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - z_j)}{[\sum_{l=1}^m y_l \phi(X_i - z_l)]^2} e_k^\top (X_i - z_j) \begin{bmatrix} \phi(X_i - z_1) \\ \vdots \\ \phi(X_i - z_m) \end{bmatrix}, \end{aligned}$$

and as a result

$$\left\| \nabla_y f_k^{(j)}(y, z) \right\|_2 \leq \frac{\sqrt{m} \|\phi\|_\infty^2}{\delta^2} \max_{i \in [N], j \in [m]} \|X_i - z_j\|_2 \leq \frac{3\sqrt{m}R}{\delta^2 (2\pi)^d}.$$

In addition, we can compute

$$\begin{aligned} \nabla_y g^{(j)}(y, z) &= -\nabla_y \left[\left(1 - \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - z_j)}{\sum_{l=1}^m y_l \phi(X_i - z_l)} \right) y_j \right] \\ &= -\left[1 - \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - z_j)}{\sum_{l=1}^m y_l \phi(X_i - z_l)} \right] e_j + \frac{1}{N} \sum_{i=1}^N \frac{y_j \phi(X_i - z_j)}{[\sum_{l=1}^m y_l \phi(X_i - z_l)]^2} \begin{bmatrix} \phi(X_i - z_1) \\ \vdots \\ \phi(X_i - z_m) \end{bmatrix} \end{aligned}$$

and for each $l \in [m]$

$$\begin{aligned} \nabla_{z_l} g^{(j)}(y, z) &= -\nabla_{z_l} \left[\left(1 - \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - z_j)}{\sum_{l=1}^m y_l \phi(X_i - z_l)} \right) y_j \right] \\ &= y_j \frac{1}{N} \sum_{i=1}^N \left[\frac{\nabla_{z_l} \phi(X_i - z_j)}{\sum_{l=1}^m y_l \phi(X_i - z_l)} - \frac{\phi(X_i - z_j) y_l \nabla_{z_l} \phi(X_i - z_l)}{[\sum_{l=1}^m y_l \phi(X_i - z_l)]^2} \right] \\ &= \frac{1}{N} \sum_{i=1}^N \left[\mathbb{1}\{l = j\} \frac{y_j \phi(X_i - z_j)}{\sum_{l=1}^m y_l \phi(X_i - z_l)} (X_i - z_j) - \frac{y_j \phi(X_i - z_j) y_l \phi(X_i - z_l)}{[\sum_{l=1}^m y_l \phi(X_i - z_l)]^2} (X_i - z_l) \right]. \end{aligned}$$

As a result, we have

$$\begin{aligned} \left\| \nabla_y g^{(j)}(y, z) \right\|_2 &\leq \left| 1 - \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - z_j)}{\sum_{l=1}^m y_l \phi(X_i - z_l)} \right| + \frac{1}{N} \sum_{i=1}^N \left| \frac{y_j \phi(X_i - z_j)}{[\sum_{l=1}^m y_l \phi(X_i - z_l)]^2} \right| \sqrt{m} \|\phi\|_\infty \\ &\leq 1 + \frac{\|\phi\|_\infty}{\delta} + \frac{(1+\varepsilon)\sqrt{m}}{\delta^2} \|\phi\|_\infty^2 = 1 + \frac{1}{(2\pi)^{d/2} \delta} + \frac{(1+\varepsilon)\sqrt{m}}{(2\pi)^d \delta^2} \end{aligned}$$

and

$$\begin{aligned} \left\| \nabla_{z_l} g^{(j)}(y, z) \right\|_2 &\leq \frac{1}{N} \sum_{i=1}^N \frac{(1+\varepsilon)\|\phi\|_\infty}{\delta} \|X_i - z_j\|_2 + \frac{1}{N} \sum_{i=1}^N \frac{(1+\varepsilon)^2 \|\phi\|_\infty^2}{\delta^2} \|X_i - z_l\|_2 \\ &\leq \frac{3(1+\varepsilon)R}{(2\pi)^{d/2} \delta} + \frac{3(1+\varepsilon)^2 R \|\phi\|_\infty^2}{(2\pi)^d \delta^2}. \end{aligned}$$

Therefore for any $k \in [d]$, $j \in [m]$, $z \in \mathbb{R}^{md}$ satisfying $\max_{j \in [m]} \|z_j\|_2 \leq 2R$ and $y \in \mathbb{R}^m$ such that $\|y - \omega_0\|_2 \leq \varepsilon$, we have

$$\begin{aligned} \left\| \nabla f_k^{(j)}(y, z) \right\|_2 &= \sqrt{\left\| \nabla_z f_k^{(j)}(y, z) \right\|_2^2 + \left\| \nabla_y f_k^{(j)}(y, z) \right\|_2^2} \\ &\leq \sqrt{(m+1) \left(\frac{9m^{-1}}{\delta^2 (2\pi)^d} R^2 \right)^2 + \frac{2}{\delta^2 (2\pi)^d} + \left(\frac{3\sqrt{m}R}{\delta^2 (2\pi)^d} \right)^2} \triangleq C_{\text{lip},f}, \end{aligned}$$

which suggests that $f^{(j)}(y, z)$ is $\sqrt{d}C_{\text{lip},f}$ -Lipschitz continuous, and

$$\begin{aligned} \left\| \nabla g^{(j)}(y, z) \right\|_2 &= \sqrt{\left\| \nabla_y g^{(j)}(y, z) \right\|_2^2 + \sum_{l=1}^m \left\| \nabla_{z_l} g^{(j)}(y, z) \right\|_2^2} \\ &\leq \sqrt{\left(1 + \frac{1}{(2\pi)^{d/2} \delta} + \frac{(1+\varepsilon)\sqrt{m}}{(2\pi)^d \delta^2} \right)^2 + m \left(\frac{3(1+\varepsilon)R}{(2\pi)^{d/2} \delta} + \frac{3(1+\varepsilon)^2 R \|\phi\|_\infty^2}{(2\pi)^d \delta^2} \right)^2} \triangleq C_{\text{lip},g}. \end{aligned}$$

which suggests that $g^{(j)}(y, z)$ is $C_{\text{lip},g}$ -Lipschitz continuous. This allows us to conclude that $h(y, z)$ is C_{Lip} -continuous in $\{(y, z) : \|y - \omega_0\|_2 \leq \varepsilon, \max_{j \in [m]} \|z_j\|_2 \leq 2R\}$, where $C_{\text{Lip}} = \sqrt{mC_{\text{lip},g}^2 + m d C_{\text{lip},f}^2}$. In addition, it is easy to check that for any $z \in \mathbb{R}^{md}$ satisfying $\max_{j \in [m]} \|z_j\|_2 \leq 2R$ and any $y \in \mathbb{R}^m$ satisfying $\|y - \omega_0\|_2 \leq \varepsilon$,

$$\max_{j \in [m]} \left\| f^{(j)}(y, z) \right\|_2 \leq \frac{3R}{\delta (2\pi)^{d/2}}, \quad \max_{j \in [m]} \left| g^{(j)}(y, z) \right| \leq \left(1 + \frac{1}{\delta (2\pi)^{d/2}} \right) (1 + \varepsilon)$$

and therefore

$$\|h(y, z)\|_2 \leq \sqrt{m \left(\frac{3R}{\delta (2\pi)^{d/2}} \right)^2 + m \left[\left(1 + \frac{1}{\delta (2\pi)^{d/2}} \right) (1 + \varepsilon) \right]^2} \triangleq M.$$

Recall that $\mu_0^{(1)}, \dots, \mu_0^{(m)}$ are i.i.d. sampled from $\text{Uniform}(\{X_i\}_{1 \leq i \leq N})$, therefore

$$(\omega_0, \mu_0) \in (y, z) : \left\{ (y, z) : \|y - \omega_0\|_2 \leq \varepsilon, \max_{j \in [m]} \|z_j\|_2 \leq 2R \right\}$$

where $\mu_0 = [\mu_0^{(j)}]_{1 \leq j \leq m}$. We are ready to apply Picard-Lindelöf theorem to show that there exists $t_0 > 0$, only depending on C_{Lip} and M , such that the ODE has a unique solution on the time interval $[0, t_0]$. We can use the same argument in the proof of Theorem 5 in Appendix E.1 to show that $\omega_t \in \Delta^{m-1}$ for all $t \in [0, t_0]$, and can use the same argument in the proof of Theorem 6 in Appendix E.2 to show that $\max_{j \in [m]} \|\mu_t^{(j)}\|_2 \leq R$ for all $t \in [0, t_0]$. Then we can repeat the same analysis as above (notice that the above proof only requires $\omega_0 \in \Delta^{m-1}$ and $\max_{j \in [m]} \|\mu_0^{(j)}\|_2 \leq R$, and t_0 only depends on universal constants C_{Lip} and M) to show that the ODE has a unique solution on $[t_0, 2t_0]$, $[2t_0, 3t_0]$, and so on. This shows the existence and uniqueness of the solution to the ODE system (3.8).

Finally we check that $(\rho_t)_{t \geq 0}$ defined as $\rho_t := \sum_{l=1}^m \omega_t^{(l)} \delta_{\mu_t^{(l)}}$ solves (4.1). Note that ρ_t is a probability measure since we have shown that $\omega_t \in \Delta^{m-1}$ for any $t \geq 0$. For any test function $\varphi(x) \in C_c^\infty$, we have

$$\begin{aligned} \frac{d}{dt} \int_{\mathbb{R}^d} \varphi(x) \rho_t(dx) &= \frac{d}{dt} \left[\sum_{j=1}^m \omega_t^{(j)} \varphi(\mu_t^{(j)}) \right] = \sum_{j=1}^m \left[\dot{\omega}_t^{(j)} \varphi(\mu_t^{(j)}) + \omega_t^{(j)} \left\langle \nabla \varphi(\mu_t^{(j)}), \dot{\mu}_t^{(j)} \right\rangle \right] \\ &= - \sum_{j=1}^m \left[1 + \delta \ell_N(\rho_t)(\mu_t^{(j)}) \right] \omega_t^{(j)} \varphi(\mu_t^{(j)}) + \sum_{j=1}^m \omega_t^{(j)} \left\langle \nabla \varphi(\mu_t^{(j)}), -\nabla \delta \ell_N(\rho_t)(\mu_t^{(j)}) \right\rangle \end{aligned}$$

$$= - \int_{\mathbb{R}^d} [1 + \delta \ell_N(\rho_t)(x)] \varphi(x) \rho_t(dx) - \int_{\mathbb{R}^d} \langle \nabla \varphi(x), \nabla \delta \ell_N(\rho_t) \rangle \rho_t(dx).$$

This proves that

$$\partial_t \rho_t = \operatorname{div}(\rho_t \nabla \delta \ell(\rho_t)) - [\delta \ell(\rho_t) + 1] \rho_t$$

holds in the sense of distributions.

F Properties of Wasserstein gradient flow

In this section, we present some preliminary results on Wasserstein gradient flow for learning Gaussian mixtures. We also discuss the implications of these results, as well as the technical difficulty of obtaining more general results.

We first establish the connection between the Wasserstein gradient flow and the classical gradient flow in the Euclidean space. Suppose we fit the data $\{X_i\}_{1 \leq i \leq N}$ using a m -component Gaussian mixture model

$$\frac{1}{m} \sum_{j=1}^m \mathcal{N}(\mu^{(j)}, I_d),$$

where $\{\mu^{(j)}\}_{1 \leq j \leq m}$ is the location of the m Gaussian components. The negative likelihood function is

$$\ell_{N,m}(\mu^{(1)}, \dots, \mu^{(m)}) := -\frac{1}{N} \sum_{i=1}^N \log \left[\frac{1}{m} \sum_{j=1}^m \phi(X_i - \mu^{(j)}) \right]. \quad (\text{F.1})$$

The gradient flow for minimizing (F.1), denoted by $(\mu_t)_{t \geq 0}$ where $\mu_t = [\mu_t^{(j)}]_{1 \leq j \leq m}$, is given by the following ODE system

$$\dot{\mu}_t^{(j)} = -\nabla_{\mu^{(j)}} \ell_{N,m}(\mu_t^{(1)}, \dots, \mu_t^{(m)}) \quad (\text{F.2})$$

with initialization $\mu_0^{(1)}, \dots, \mu_0^{(m)} \stackrel{i.i.d.}{\sim} \text{Uniform}(\{X_i\}_{1 \leq i \leq N})$. The following theorem shows that the gradient flow (F.2) captures the evolution of the location of particles in the Wasserstein gradient flow (4.8) initialized from a discrete distribution $\frac{1}{m} \sum_{l=1}^m \delta_{\mu_0^{(l)}}$. The proof is deferred to Appendix F.1.

Theorem 1. *Consider the Euclidean gradient flow $(\mu_t)_{t \geq 0}$ in (F.2). Then the flow $(\rho_t)_{t \geq 0}$ defined as*

$$\rho_t := \frac{1}{m} \sum_{l=1}^m \delta_{\mu_t^{(l)}} \quad (\text{F.3})$$

is the Wasserstein gradient flow, i.e. (F.3) is a distributional solution to the PDE (4.8).

Similar connection can also be established for the gradient descent algorithm for minimizing (F.1) and the particle Wasserstein gradient descent (cf. Algorithm 3), which is omitted for brevity.

Then we focus on the infinite sample limit of Wasserstein gradient flow and analyze its convergence property. The population level loss function is

$$\ell_\infty(\rho) = -\mathbb{E}_{X \sim \rho^* * \phi} \{\log[\rho * \phi(X)]\} = \text{KL}(\rho^* * \phi \| \rho * \phi) + \text{const}. \quad (\text{F.4})$$

In Appendix C.1 we have computed that

$$\delta \ell_\infty(\rho) = - \int_{\mathbb{R}^d} \frac{\rho^* * \phi(y)}{\rho * \phi(y)} \phi(x - y) dy. \quad (\text{F.5})$$

We know that the Wasserstein gradient flow $(\rho_t)_{t \geq 0}$ with respect to $\ell_\infty(\rho)$ is described by the following PDE:

$$\partial_t \rho_t = \operatorname{div}(\rho_t \nabla \delta \ell_\infty(\rho_t)) \quad (\text{F.6})$$

with $\rho_0 = \rho^* * \mathcal{N}(0, I_d)$, which is the data distribution when we have infinite samples. This Wasserstein gradient flow has the following particle interpretation: suppose at time $t = 0$ we initialize a particle $x_0 \sim \rho_0$ in the vector field $(v_t)_{t \geq 0}$ where $v_t = -\nabla \delta \ell_\infty(\rho_t)$, namely

$$\dot{x}_t = v_t(x_t),$$

then $x_t \sim \rho_t$, namely the marginal distribution of $(x_t)_{t \geq 0}$ evolves according to the Wasserstein gradient flow.

The following theorem shows that, when the true mixing distribution ρ^* is a singleton (we assume without loss of generality that $\rho^* = \delta_0$), Wasserstein gradient flow converges to ρ^* . The proof can be found in Appendix F.2.

Theorem 2. *Consider the Wasserstein gradient flow in (F.6) with $\rho^* = \delta_0$. For any $\varepsilon < 1$, we have $\int_{\mathbb{R}^d} \|x\|_2^2 \rho_t(dx) = O(\varepsilon)$ as long as*

$$t \geq \exp(2d) \varepsilon^{-1 - \max\{8, \sqrt{8d}\}}.$$

Although Theorem 2 only focuses on the case when ρ^* is a singleton, the convergence result already provides some intuition about the behavior of Wasserstein gradient flow in more general setting. Consider a well-separated Gaussian mixture model with K components. Assume that the mixing distribution is $\rho^* = \sum_{j=1}^K \omega_j^* \delta_{\mu_j^*}$, and the location of each Gaussian components, $\{\mu_j^*\}_{1 \leq j \leq K}$, are well-separated. Since the push-forward mapping $v_t = -\nabla \delta \ell_\infty(\rho_t)$ is localized (see F.5), there exists some $T > 0$ such that the Wasserstein gradient flow (F.6) initialized from $\rho^* * \mathcal{N}(0, I_d)$ can be approximated, up to time T , by

$$\rho_t \approx \sum_{j=1}^K \omega_j^* \rho_t^{(j)} \quad \forall t \in [0, T],$$

where for each $j \in [K]$, $\rho_t^{(j)}$ is the Wasserstein gradient flow $\partial_t \rho_t^{(j)} = \text{div}(\rho_t^{(j)} \nabla \delta \ell_\infty(\rho_t^{(j)}))$ with initialization $\rho_0^{(j)} = \mathcal{N}(\mu_j^*, I_d)$. This suggests that Wasserstein gradient flow approximately converges to ρ^* since, by Theorem 2, each $\rho_t^{(j)}$ converges to $\delta_{\mu_j^*}$. However this observation also suggests that Wasserstein gradient flow is not robust to weight mismatch. Consider initializing the Wasserstein gradient flow (F.6) with $\rho_0 = \tilde{\rho} * \mathcal{N}(0, 1)$, where $\tilde{\rho} = \sum_{j=1}^K \tilde{\omega}_j \delta_{\mu_j^*}$ is a mixing distribution with correct support $\{\mu_j^*\}_{1 \leq j \leq K}$ but wrong weights $\{\tilde{\omega}_j\}_{1 \leq j \leq K} \neq \{\omega_j^*\}_{1 \leq j \leq K}$. Then we also have

$$\rho_t \approx \sum_{j=1}^K \tilde{\omega}_j \rho_t^{(j)} \quad \forall t \in [0, T],$$

which shows that ρ_t approximately converges to $\tilde{\rho}$ instead of ρ^* when $0 \leq t \leq T$. Note that the time length T that such approximations are valid can be arbitrarily large as long as the separation $\min_{i \neq j} \|\mu_i^* - \mu_j^*\|_2 \rightarrow \infty$. The above discussion suggests that using the correct initial weights are important for Wasserstein gradient flow to converge to the true mixing distribution.

We also would like to compare the convergence rate in Theorem 2 to a benchmark provided by the Bures-Wasserstein gradient flow. The Bures-Wasserstein gradient flow is defined on the space of non-degenerate Gaussian distributions on \mathbb{R}^d , denoted by $\text{BW}(\mathbb{R}^d) = \mathbb{R}^d \times \mathbb{S}_{++}^d$ (where we identify a non-degenerate Gaussian distribution $\nu = \mathcal{N}(\mu, \Sigma)$ with $(\mu, \Sigma) \in \mathbb{R}^d \times \mathbb{S}_{++}^d$) equipped with the Wasserstein distance (3.3), which has the following closed form expression

$$d_W^2(\nu_1, \nu_2) = \|\mu_1 - \mu_2\|_2^2 + \text{tr} \left[\Sigma_1 + \Sigma_2 - 2 \left(\Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2} \right)^{1/2} \right]$$

when $\nu_1 = \mathcal{N}(\mu_1, \Sigma_1)$ and $\nu_2 = \mathcal{N}(\mu_2, \Sigma_2)$ are both non-degenerate Gaussians. The Bures-Wasserstein gradient flow $(\nu_t)_{t \geq 0}$ can be viewed as the Wasserstein gradient flow $(\rho_t)_{t \geq 0}$ constrained to lie on $\text{BW}(\mathbb{R}^d)$. We refer interested readers to Altschuler et al. (2021); Lambert et al. (2022) for more detailed discussion. We can see from the proof of Theorem 2 that the push-forward mapping $v_t(x)$ of Wasserstein gradient flow decays exponentially fast as $\|x\|_2 \rightarrow \infty$, this will make the Wasserstein gradient flow $(\rho_t)_{t \geq 0}$ becomes more and more heavy-tailed. However the push-forward mapping of Bures-Wasserstein gradient flow is always

linear, and the Bures-Wasserstein gradient flow $(\nu_t)_{t \geq 0}$ is always Gaussian. For example, in Appendix F.3 we can compute the push forward mapping explicitly for the two gradient flows at $t = 0$:

$$v_0(x) = -\frac{1}{3} \left(\frac{4}{3}\right)^{d/2} \exp\left(-\frac{\|x\|_2^2}{6}\right) x \quad (\text{Wasserstein}),$$

$$v_0(x) = \frac{x}{4} \quad (\text{Bures-Wasserstein}).$$

Therefore it is natural to expect that Bures-Wasserstein gradient flow $(\nu_t)_{t \geq 0}$ initialized from $\nu_0 = \mathcal{N}(0, I_d)$ converges faster than Wasserstein gradient flow $(\rho_t)_{t \geq 0}$ initialized from $\rho_0 = \delta_0$. In Appendix F.3 we show that the Bures-Wasserstein gradient flow $(\nu_t = \mathcal{N}(\mu_t, \Sigma_t))_{t \geq 0}$ is characterized by the following ODE:

$$\begin{aligned} \mu_t &= 0 \\ \dot{\Sigma}_t &= -2(\Sigma_t + I_d)^{-1} \Sigma_t^2 (\Sigma_t + I_d)^{-1}. \end{aligned}$$

We also show that Σ_t is sandwiched between

$$\frac{1}{1+2t} I \preceq \Sigma_t \preceq \frac{2}{2+t} I,$$

and as a result $\int_{\mathbb{R}^d} \|x\|_2^2 \nu_t(dx) = O(d/t)$. Since Bures-Wasserstein gradient flow is not converging exponentially fast (we can see that the convergence rate is polynomial in t), we conjecture that Wasserstein gradient flow does not enjoy exponential convergence as well.

Lastly, we numerically show in Figure 1 that the loss function $\ell_\infty(\rho)$ (cf. (F.4)) is not geodesically convex (Ambrosio et al., 2008) even when $\rho^* = \delta_0$. We can also check that Polyak-Łojasiewicz (PL) inequality

$$\forall \rho: \quad \|\nabla_W \ell_\infty(\rho)\|_\rho^2 \geq C_{\text{PL}} [\ell_\infty(\rho) - \ell_\infty(\rho^*)] \quad \text{for some } C_{\text{PL}} > 0$$

does not hold in general: consider $\rho^* = \frac{1}{2}\delta_{-1} + \frac{1}{2}\delta_1$ and $\rho = \delta_0$, then it is straightforward to check that $\nabla_W \ell_\infty(\rho) = 0$ but $\ell_\infty(\rho) > \ell_\infty(\rho^*)$. Therefore we cannot use standard proof technique (e.g. Ambrosio et al. (2008) when the loss function is geodesically convex, or Chewi et al. (2020) when there is a PL inequality) to show exponential convergence for the Wasserstein gradient flow (F.6).

F.1 Proof of Theorem 1

It is straightforward to compute the gradient of $\ell_{N,m}$. For any $j \in [m]$, we have

$$\begin{aligned} \nabla_{\mu^{(j)}} \ell_{N,m}(\mu^{(1)}, \dots, \mu^{(m)}) &= -\frac{1}{N} \sum_{i=1}^N \frac{1}{\sum_{l=1}^m \phi(X_i - \mu^{(l)})} \nabla_{\mu^{(j)}} \phi(X_i - \mu^{(j)}) \\ &= -\frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - \mu^{(l)})}{\sum_{l=1}^m \phi(X_i - \mu^{(l)})} (X_i - \mu^{(l)}). \end{aligned}$$

Therefore the Euclidean gradient flow (F.2) is given by

$$\dot{\mu}_t^{(j)} = -\nabla_{\mu^{(j)}} \ell_{N,m}(\mu_t^{(1)}, \dots, \mu_t^{(m)}) = \frac{1}{N} \sum_{i=1}^N \frac{\phi(X_i - \mu_t^{(l)})}{\sum_{l=1}^m \phi(X_i - \mu_t^{(l)})} (X_i - \mu_t^{(l)}).$$

Then we can invoke Theorem 6 to finish the proof.

F.2 Proof of Theorem 2

Step 1: characterizing the push-forward mapping. First of all, it is straightforward to check that $(\rho_t)_{t \geq 0}$ is spherically symmetric for all $t \geq 0$, namely $\rho_t(dx)$ only depends on $\|x\|_2$. The push-forward mapping $v_t(x) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ at time t is

$$v_t(x) = -\nabla \delta \ell_\infty(\rho_t)(x) = \nabla_x \int_{\mathbb{R}^d} \frac{\rho^* * \phi(y)}{\rho_t * \phi(y)} \phi(x - y) dy$$

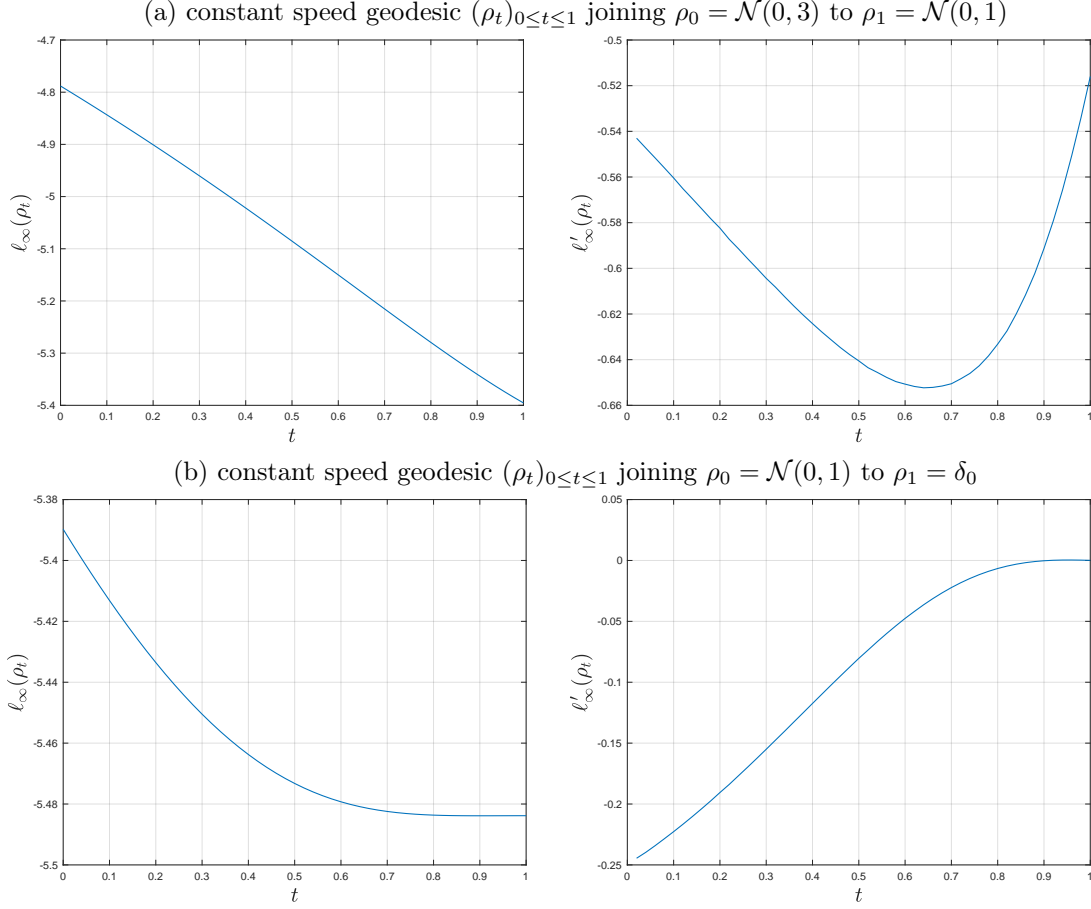


Figure 1: The loss function $\ell_\infty(\rho_t)$ or its derivative $\ell'_\infty(\rho_t)$ vs. t . In Figures (a), $(\rho_t)_{0 \leq t \leq 1}$ is the constant speed geodesic joining $\rho_0 = \mathcal{N}(0, 3)$ to $\rho_1 = \mathcal{N}(0, 1)$. In Figures (b), $(\rho_t)_{0 \leq t \leq 1}$ is the constant speed geodesic joining $\rho_0 = \mathcal{N}(0, 1)$ to $\rho_1 = \delta_0$. This shows that $\ell_\infty(\rho)$ is not globally geodesically convex, but might be locally geodesically convex around $\rho^* = \delta_0$.

$$\begin{aligned}
&= - \int_{\mathbb{R}^d} \frac{\rho^\star * \phi(y)}{\rho_t * \phi(y)} (x - y) \phi(x - y) dy \\
&= \int_{\mathbb{R}^d} \nabla_y \frac{\rho^\star * \phi(y)}{\rho_t * \phi(y)} \phi(y - x) dy \\
&= \int_{\mathbb{R}^d} h_t(y) \phi(y - x) dy,
\end{aligned} \tag{F.7}$$

where the penultimate line follows from Stein's lemma or Gaussian integration by parts, and $h_t(y)$ in the last line is defined as

$$\begin{aligned}
h_t(y) &:= \nabla \frac{\phi(y)}{\rho_t * \phi(y)} = \nabla_y \frac{1}{\int \exp\left(-\frac{1}{2}\|z\|_2^2 + y^\top z\right) \rho_t(dz)} \\
&= - \frac{\int \exp\left(-\frac{1}{2}\|z\|_2^2 + y^\top z\right) z \rho_t(dz)}{\left[\int \exp\left(-\frac{1}{2}\|z\|_2^2 + y^\top z\right) \rho_t(dz)\right]^2} = -\phi(y) \frac{\int \phi(y - z) z \rho_t(dz)}{\left[\int \phi(y - z) \rho_t(dz)\right]^2} \\
&= -\phi(y) \frac{\int \phi(y - z) z \rho_t(dz)}{[\rho_t * \phi(y)]^2} = -\frac{\phi(y)}{\rho_t * \phi(y)} \cdot \frac{\int \phi(y - z) z \rho_t(dz)}{\rho_t * \phi(y)}.
\end{aligned} \tag{F.8}$$

For any $y \in \mathbb{R}^d$, we can compute

$$\begin{aligned}
\int \phi(y-z) z \rho_t(dz) &= \int_{y^\top z > 0} \phi(y-z) z \rho_t(dz) + \int_{y^\top z < 0} \phi(y-z) z \rho_t(dz) + \int_{y^\top z = 0} \phi(y-z) z \rho_t(dz) \\
&\stackrel{(i)}{=} \int_{y^\top z > 0} [\phi(y-z) - \phi(y+z)] z \rho_t(dz) \\
&= \phi(y) \int_{y^\top z > 0} [\exp(y^\top z) - \exp(-y^\top z)] z \exp\left(-\frac{1}{2} \|z\|_2^2\right) \rho_t(dz) \\
&\stackrel{(ii)}{=} \phi(y) \int_{y^\top z > 0} [\exp(y^\top z) - \exp(-y^\top z)] \frac{y^\top z}{\|y\|_2^2} y \exp\left(-\frac{1}{2} \|z\|_2^2\right) \rho_t(dz) \\
&= \underbrace{\int_{y^\top z > 0} [\exp(y^\top z) - \exp(-y^\top z)] \frac{y^\top z}{\|y\|_2^2} \exp\left(-\frac{1}{2} \|z\|_2^2\right) \rho_t(dz)}_{=: a_t} \cdot \phi(y) y. \tag{F.9}
\end{aligned}$$

Here (i) and (ii) both follow from the spherical symmetry of ρ_t , and it is straightforward to check that the integral in the last line does not depend on y due to the spherical symmetry of ρ_t , therefore a_t is a universal constant that is independent of y . Note that when $y^\top z > 0$, we have $\exp(y^\top z) - \exp(-y^\top z) \geq 2y^\top z$, therefore

$$a_t \geq 2 \int_{y^\top z > 0} \frac{(y^\top z)^2}{\|y\|_2^2} \exp\left(-\frac{1}{2} \|z\|_2^2\right) \rho_t(dz) = \int_{\mathbb{R}^d} \frac{(y^\top z)^2}{\|y\|_2^2} \exp\left(-\frac{1}{2} \|z\|_2^2\right) \rho_t(dz).$$

Since a_t does not depend on y , we take $y = e_i$ for $i \in [d]$ to achieve

$$a_t \geq \int_{\mathbb{R}^d} z_i^2 \exp\left(-\frac{1}{2} \|z\|_2^2\right) \rho_t(dz), \quad \forall i \in [d].$$

By taking average over d , we have

$$a_t \geq \frac{1}{d} \int_{\mathbb{R}^d} \|z\|_2^2 \exp\left(-\frac{1}{2} \|z\|_2^2\right) \rho_t(dz) = \frac{m_t}{d}, \tag{F.10}$$

where we define

$$m_t := \int \|z\|_2^2 \exp\left(-\frac{1}{2} \|z\|_2^2\right) \rho_t(dz).$$

Taking (F.8), (F.9) and (F.10) collectively gives

$$h_t(y) = -a_t y \left[\frac{\phi(y)}{\rho_t * \phi(y)} \right]^2, \quad \text{where} \quad a_t \geq \frac{m_t}{d}.$$

Then we use (F.7) to characterize the push-forward mapping:

$$\begin{aligned}
v_t(x) &= \int_{y^\top x > 0} h(y) \phi(y-x) dy + \int_{y^\top x < 0} h(y) \phi(y-x) dy + \int_{y^\top x = 0} h(y) \phi(y-x) dy \\
&\stackrel{(i)}{=} \int_{y^\top x > 0} h(y) [\phi(y-x) - \phi(-y-x)] dy \\
&= -\phi(x) a_t \int_{y^\top x > 0} y \left[\frac{\phi(y)}{\rho_t * \phi(y)} \right]^2 [\exp(y^\top x) - \exp(-y^\top x)] \exp\left(-\frac{1}{2} \|y\|_2^2\right) dy \\
&\stackrel{(ii)}{=} -\phi(x) a_t \int_{y^\top x > 0} \frac{x^\top y}{\|x\|_2^2} x \left[\frac{\phi(y)}{\rho_t * \phi(y)} \right]^2 [\exp(y^\top x) - \exp(-y^\top x)] \exp\left(-\frac{1}{2} \|y\|_2^2\right) dy \\
&= -a_t \underbrace{\int_{y^\top x > 0} \frac{x^\top y}{\|x\|_2^2} \left[\frac{\phi(y)}{\rho_t * \phi(y)} \right]^2 [\exp(y^\top x) - \exp(-y^\top x)] \exp\left(-\frac{1}{2} \|y\|_2^2\right) dy}_{=: b_t} \cdot \phi(x) x.
\end{aligned}$$

Similar to (F.9), here (i) and (ii) both follow from the spherical symmetry of ρ_t , and the integral in the last line does not depend on x due to the spherical symmetry of ρ_t , as a result b_t is a universal constant that is independent of x . Note that when $y^\top x > 0$, we have $\exp(y^\top x) - \exp(-y^\top x) \geq 2y^\top x$, therefore

$$\begin{aligned} b_t &\geq 2 \int_{y^\top x > 0} \frac{x^\top y}{\|x\|_2^2} \left[\frac{\phi(y)}{\rho_t * \phi(y)} \right]^2 y^\top x \exp\left(-\frac{1}{2} \|y\|_2^2\right) dy \\ &= \int_{\mathbb{R}^d} \frac{(x^\top y)^2}{\|x\|_2^2} \left[\frac{\phi(y)}{\rho_t * \phi(y)} \right]^2 \exp\left(-\frac{1}{2} \|y\|_2^2\right) dy. \end{aligned}$$

Note that $\|\rho_t * \phi\|_\infty \leq \|\phi\|_\infty \leq (2\pi)^{-d/2}$, and as a result

$$b_t \geq \int \frac{(x^\top y)^2}{\|x\|_2^2} \exp\left(-\frac{3}{2} \|y\|_2^2\right) dy = \frac{1}{3} \left(\frac{2\pi}{3}\right)^{d/2}$$

Therefore we have

$$v_t(x) = -a_t b_t \phi(x) x = -c_t \phi(x) x \quad (\text{F.11})$$

where

$$c_t := a_t b_t \geq \frac{1}{3d} \left(\frac{2\pi}{3}\right)^{d/2} m_t. \quad (\text{F.12})$$

Step 2: showing the convergence of Wasserstein gradient flow. Recall the particle interpretation of Wasserstein gradient flow as follows: let $x_0 \sim \rho_0 = \mathcal{N}(0, I_d)$ and $\dot{x}_t = v_t(x_t)$, then for any $t \geq 0$ we have $x_t \sim \rho_t$. This allows us to compute

$$\begin{aligned} \partial_t \mathbb{E} \left[\|x_t\|_2^2 \right] &= 2\mathbb{E} [\langle x_t, \dot{x}_t \rangle] = 2\mathbb{E} [\langle x_t, v_t(x_t) \rangle] \stackrel{(i)}{=} -2c_t \mathbb{E} \left[\|x_t\|_2^2 \phi(x_t) \right] \\ &\stackrel{(ii)}{\leq} -\frac{2}{3d} \left(\frac{2\pi}{3}\right)^{d/2} m_t \mathbb{E} \left[\|x_t\|_2^2 \phi(x_t) \right] \\ &\stackrel{(iii)}{=} -\frac{2}{3d} \left(\frac{4\pi^2}{3}\right)^{d/2} \mathbb{E}^2 \left[\|x_t\|_2^2 \phi(x_t) \right], \end{aligned} \quad (\text{F.13})$$

where (i) follows from (F.11), (ii) utilizes (F.12), and (iii) holds since

$$m_t = \int \|z\|_2^2 \exp\left(-\frac{1}{2} \|z\|_2^2\right) \rho_t(dz) = (2\pi)^{d/2} \mathbb{E} \left[\|x_t\|_2^2 \phi(x_t) \right].$$

For any $\tau > 0$, by Cauchy-Schwarz inequality we have

$$\mathbb{E} \left[\|x_0\|_2^2 \mathbf{1} \left\{ \|x_0\|_2^2 > d + \tau \right\} \right] \stackrel{(i)}{\leq} \left[\mathbb{E} \|x_0\|_2^4 \right]^{1/2} \left[\mathbb{P} \left(\|x_0\|_2^2 > d + \tau \right) \right]^{1/2}.$$

Note that $\|x_0\|_2^2 \sim \chi^2(d)$, therefore $\mathbb{E} \|x_0\|_2^4 = \text{var}(\|x_0\|_2^2) + (\mathbb{E} \|x_0\|_2^2)^2 = 2d + d^2$. In addition, by the tail probability bound for χ^2 random variables (e.g. [Wainwright \(2019, equation \(2.18\)\)](#)), we have

$$\mathbb{P} \left(\|x_0\|_2^2 > d + \tau \right) \leq \exp \left(-\min \left\{ \frac{\tau^2}{8d}, \frac{\tau}{8} \right\} \right).$$

Therefore we have

$$\begin{aligned} \mathbb{E} \left[\|x_0\|_2^2 \mathbf{1} \left\{ \|x_0\|_2^2 > d + \tau \right\} \right] &\leq \sqrt{2d + d^2} \exp \left(-\min \left\{ \frac{\tau^2}{8d}, \frac{\tau}{8} \right\} \right) \\ &\leq (d + 1) \exp \left(-\min \left\{ \frac{\tau^2}{8d}, \frac{\tau}{8} \right\} \right) \leq \varepsilon \end{aligned}$$

as long as we choose

$$\tau \triangleq \max \left\{ 8 \log \frac{d+1}{\varepsilon}, \sqrt{8d \log \frac{d+1}{\varepsilon}} \right\}.$$

Since the push forward mapping $v_t(x)$ is always pointing towards zero (cf. (F.11) and (F.12)), we know that $\|x_t\|_2$ is non-increasing in t . Therefore we have

$$\begin{aligned} \mathbb{E} \left[\|x_t\|_2^2 \phi(x_t) \right] &\geq \mathbb{E} \left[\|x_t\|_2^2 \phi(x_t) \mathbf{1} \left\{ \|x_0\|_2^2 \leq d + \tau \right\} \right] \\ &\stackrel{(i)}{\geq} (2\pi)^{-d/2} \exp \left(-\frac{d+\tau}{2} \right) \mathbb{E} \left[\|x_t\|_2^2 \mathbf{1} \left\{ \|x_0\|_2^2 \leq d + \tau \right\} \right] \\ &\geq (2\pi)^{-d/2} \exp \left(-\frac{d+\tau}{2} \right) \left(\mathbb{E} \left[\|x_t\|_2^2 \right] - \mathbb{E} \left[\|x_t\|_2^2 \mathbf{1} \left\{ \|x_0\|_2^2 > d + \tau \right\} \right] \right) \\ &\stackrel{(ii)}{\geq} (2\pi)^{-d/2} \exp \left(-\frac{d+\tau}{2} \right) \left(\mathbb{E} \left[\|x_t\|_2^2 \right] - \mathbb{E} \left[\|x_0\|_2^2 \mathbf{1} \left\{ \|x_0\|_2^2 > d + \tau \right\} \right] \right) \\ &\geq (2\pi)^{-d/2} \exp \left(-\frac{d+\tau}{2} \right) \left(\mathbb{E} \left[\|x_t\|_2^2 \right] - \varepsilon \right), \end{aligned} \tag{F.14}$$

where both (i) and (ii) follows from the fact that $\|x_t\|_2$ is non-increasing. Taking (F.13) and (F.14) collectively gives

$$\begin{aligned} \partial_t \mathbb{E} \left[\|x_t\|_2^2 \right] &\leq -\frac{2}{3d} \left(\frac{4\pi^2}{3} \right)^{d/2} (2\pi)^{-d} \exp[-(d+\tau)] \left(\mathbb{E} \left[\|x_t\|_2^2 \right] - \varepsilon \right)^2 \\ &= -\frac{2}{3d} \left(\frac{1}{3} \right)^{d/2} \exp[-(d+\tau)] \left(\mathbb{E} \left[\|x_t\|_2^2 \right] - \varepsilon \right)^2. \end{aligned}$$

Let $f(t) = \mathbb{E}[\|x_t\|_2^2]$, we know that $f(0) = d$ and

$$\frac{df}{dt} \leq -\frac{2}{3d} \left(\frac{1}{3} \right)^{d/2} \exp[-(d+\tau)] (f - \varepsilon)^2.$$

Solving this ordinary differential inequality gives

$$\frac{1}{f(t) - \varepsilon} - \frac{1}{f(0) - \varepsilon} \geq \frac{2}{3d} \left(\frac{1}{3} \right)^{d/2} \exp[-(d+\tau)] t,$$

which is equivalent to

$$\mathbb{E} \left[\|x_t\|_2^2 \right] \leq \varepsilon + \left\{ \frac{2}{3d} \left(\frac{1}{3} \right)^{d/2} \exp(-d-\tau) t + \frac{1}{d-\varepsilon} \right\}^{-1}.$$

Then we immediately know that $\mathbb{E}[\|x_t\|_2^2] \leq O(\varepsilon)$ as long as

$$t \geq \exp(2d) \varepsilon^{-1-\max\{8, \sqrt{8d}\}}.$$

F.3 Calculation for Bures-Wasserstein gradient flow

Define $\ell(\mu, \Sigma) = \ell_\infty(\rho)$ where we parameterize $\rho = \mathcal{N}(\mu, \Sigma)$. Then we can compute

$$\begin{aligned} \ell(\mu, \Sigma) &= - \int \log \left[(2\pi)^{-d/2} [\det(\Sigma + I_d)]^{-1/2} \exp \left(-\frac{1}{2} (x - \mu)^\top (\Sigma + I_d)^{-1} (x - \mu) \right) \right] \phi(x) dx + \text{constant} \\ &= \frac{1}{2} \log \det(\Sigma + I_d) + \int \frac{1}{2} (x - \mu)^\top (\Sigma + I_d)^{-1} (x - \mu) \phi(x) dx + \text{constant} \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} \log \det (\Sigma + I_d) + \frac{1}{2} \mathbb{E}_{x \sim \mathcal{N}(0, I)} \left[(x - \mu)^\top (\Sigma + I_d)^{-1} (x - \mu) \right] + \text{constant} \\
&= \frac{1}{2} \log \det (\Sigma + I_d) + \frac{1}{2} \text{tr} \left[(\Sigma + I_d)^{-1} \right] + \frac{1}{2} \mu^\top (\Sigma + I_d)^{-1} \mu + \text{constant}.
\end{aligned}$$

Then we can compute the Euclidean gradient of $\ell(\mu, \Sigma)$ as follows:

$$\begin{aligned}
\nabla_\mu \ell(\mu, \Sigma) &= (\Sigma + I_d)^{-1} \mu, \\
\nabla_\Sigma \ell(\mu, \Sigma) &= \frac{1}{2} (\Sigma + I_d)^{-1} - \frac{1}{2} (\Sigma + I_d)^{-2} - \frac{1}{2} (\Sigma + I_d)^{-1} \mu \mu^\top (\Sigma + I_d)^{-1} \\
&= \frac{1}{2} (\Sigma + I_d)^{-1} (\Sigma + I_d - I_d - \mu \mu^\top) (\Sigma + I_d)^{-1} \\
&= \frac{1}{2} (\Sigma + I_d)^{-1} (\Sigma - \mu \mu^\top) (\Sigma + I_d)^{-1}.
\end{aligned}$$

According to [Lambert et al. \(2022, Appendix B.3\)](#), when initialized from $(\mu_0, \Sigma_0) = (0, I_d)$, the Bures-Wasserstein gradient flow can be described using the following ODE:

$$\begin{aligned}
\dot{\mu}_t &= -(\Sigma_t + I_d)^{-1} \mu_t \\
\dot{\Sigma}_t &= -\Sigma_t (\Sigma_t + I_d)^{-1} [\Sigma_t - \mu \mu^\top] (\Sigma_t + I_d)^{-1} - (\Sigma_t + I_d)^{-1} [\Sigma_t - \mu \mu^\top] (\Sigma_t + I_d)^{-1} \Sigma_t
\end{aligned}$$

with initial condition $\mu_0 = 0$ and $\Sigma_0 = I_d$. It is straightforward to check that $\mu_t = 0$ for all $t \geq 0$, and the dynamic of Σ_t is governed by

$$\begin{aligned}
\dot{\Sigma}_t &= -\Sigma_t (\Sigma_t + I_d)^{-1} \Sigma_t (\Sigma_t + I_d)^{-1} - (\Sigma_t + I_d)^{-1} \Sigma_t (\Sigma_t + I_d)^{-1} \Sigma_t \\
&= -\Sigma_t (\Sigma_t + I_d)^{-1} + 2 (\Sigma_t + I_d)^{-1} \Sigma_t (\Sigma_t + I_d)^{-1} - (\Sigma_t + I_d)^{-1} \Sigma_t \\
&= -2I_d + 2 (\Sigma_t + I_d)^{-1} + 2 (\Sigma_t + I_d)^{-1} \Sigma_t (\Sigma_t + I_d)^{-1} \\
&= -2 (\Sigma_t + I_d)^{-1} \Sigma_t^2 (\Sigma_t + I_d)^{-1}
\end{aligned}$$

with initial condition $\Sigma_0 = I_d$. We can check that the off-diagonal entries of Σ_t are always zero, and its diagonal entries are identical and evolves according to the following ODE

$$\dot{\sigma}_t = -2 \frac{\sigma_t^2}{(\sigma_t + 1)^2}$$

with initial condition $\sigma_0 = 1$. It is straightforward to check that σ_t is monotonically decreasing and is always non-negative, namely $0 \leq \sigma_t \leq 1$ always holds. Therefore we have

$$-2\sigma_t^2 \leq \dot{\sigma}_t \leq -\frac{1}{2}\sigma_t^2.$$

This gives

$$\frac{1}{1+2t} \leq \sigma_t \leq \frac{2}{2+t},$$

and therefore

$$\frac{1}{1+2t} I \preceq \Sigma_t \preceq \frac{2}{2+t} I,$$

which suggests that ρ_t converges to ρ^* at the speed of $O(d/t)$.

When $t = 0$, we can compute the push forward mapping of Wasserstein gradient flow explicitly, which intuitively explains why Wasserstein gradient flow does not converge exponentially fast. We first compute

$$\nabla \frac{\rho^* * \phi(y)}{\rho_0 * \phi(y)} = \nabla \frac{(\det I)^{-d/2} \exp\left(-\frac{1}{2} \|y\|_2^2\right)}{(\det 2I)^{-d/2} \exp\left(-\frac{1}{4} \|y\|_2^2\right)} = 2^{d/2} \nabla \exp\left(-\frac{1}{4} \|y\|_2^2\right) = -2^{d/2-1} y \exp\left(-\frac{1}{4} \|y\|_2^2\right),$$

then the push forward mapping at $t = 0$ is given by $x \mapsto v_0(x)$ where

$$\begin{aligned}
v_0(x) &= \int \left(\nabla_y \frac{\rho^\star * \phi(y)}{\rho_0 * \phi(y)} \right) \phi(y - x) dy = -2^{d/2-1} \int y \exp\left(-\frac{1}{4} \|y\|_2^2\right) \cdot \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2} \|x - y\|_2^2\right) dy \\
&= -\frac{2^{d/2-1}}{(2\pi)^{d/2}} \int y \exp\left(-\frac{1}{2} \|x\|_2^2 + x^\top y - \frac{3}{4} \|y\|_2^2\right) dy = -\frac{2^{d/2-1}}{(2\pi)^{d/2}} \int y \exp\left(-\frac{1}{6} \|x\|_2^2 - \frac{3}{4} \left\|y - \frac{2}{3}x\right\|_2^2\right) dy \\
&= -2^{d/2-1} \left(\frac{2}{3}\right)^{d/2} \exp\left(-\frac{1}{6} \|x\|_2^2\right) \int \frac{1}{(2\pi)^{d/2} (2/3)^{d/2}} y \exp\left(-\frac{3}{4} \left\|y - \frac{2}{3}x\right\|_2^2\right) dy \\
&= -\frac{1}{3} \left(\frac{4}{3}\right)^{d/2} \exp\left(-\frac{1}{6} \|x\|_2^2\right) x.
\end{aligned}$$

On the other hand, in view of [Lambert et al. \(2022, Appendix B.3\)](#), the Bures-Wasserstein gradient at time $t = 0$ is given by

$$\nabla_{\text{BW}} \ell_\infty(\rho_0) = \begin{bmatrix} \nabla_\mu \ell(\mu_0, \Sigma_0) \\ 2\nabla_\Sigma \ell(\mu_0, \Sigma_0) \end{bmatrix} = \begin{bmatrix} 0 \\ \frac{1}{4} I_d \end{bmatrix},$$

and therefore the push forward mapping at $t = 0$ is given by $x \mapsto x/4$.

References

- Altschuler, J., Chewi, S., Gerber, P. R., and Stromme, A. (2021). Averaging on the bures-wasserstein manifold: dimension-free convergence of gradient descent. *Advances in Neural Information Processing Systems*, 34:22132–22145.
- Ambrosio, L., Gigli, N., and Savaré, G. (2008). *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media.
- Bubeck, S. (2015). *Convex optimization: algorithms and complexity*. Now Publishers Inc.
- Chewi, S., Maunu, T., Rigollet, P., and Stromme, A. J. (2020). Gradient descent algorithms for bures-wasserstein barycenters. In *Conference on Learning Theory*, pages 1276–1304. PMLR.
- Gallouët, T. O. and Monsaingeon, L. (2017). A jko splitting scheme for kantorovich–fisher–rao gradient flows. *SIAM Journal on Mathematical Analysis*, 49(2):1100–1130.
- Lambert, M., Chewi, S., Bach, F., Bonnabel, S., and Rigollet, P. (2022). Variational inference via wasserstein gradient flows. *arXiv preprint arXiv:2205.15902*.
- Lu, Y., Lu, J., and Nolen, J. (2019). Accelerating langevin sampling with birth-death. *arXiv preprint arXiv:1905.09863*.
- Salim, A., Korba, A., and Luise, G. (2020). The wasserstein proximal gradient algorithm. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems*, volume 33, pages 12356–12366. Curran Associates, Inc.
- Wainwright, M. J. (2019). *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press.