

LINEAR CONVERGENCE OF PROXIMAL DESCENT SCHEMES ON THE WASSERSTEIN SPACE

RAZVAN-ANDREI LASCU, MATEUSZ B. MAJKA, DAVID ŠIŠKA, AND ŁUKASZ SZPRUCH

ABSTRACT. We investigate proximal descent methods, inspired by the minimizing movement scheme introduced by Jordan, Kinderlehrer and Otto, for optimizing entropy-regularized functionals on the Wasserstein space. We establish linear convergence under flat convexity assumptions, thereby relaxing the common reliance on geodesic convexity. Our analysis circumvents the need for discrete-time adaptations of the Evolution Variational Inequality (EVI). Instead, we leverage a uniform logarithmic Sobolev inequality (LSI) and the entropy “sandwich” lemma, extending the analysis from [27, 12]. The major challenge in the proof via LSI is to show that the relative Fisher information $I(\cdot|\pi)$ is well-defined at every step of the scheme. Since the relative entropy is not Wasserstein differentiable, we prove that along the scheme the iterates belong to a certain class of Sobolev regularity, and hence the relative entropy $\text{KL}(\cdot|\pi)$ has a unique Wasserstein sub-gradient, and that the relative Fisher information is indeed finite.

1. INTRODUCTION

We consider the problem of minimizing an entropy-regularized flat-convex function

$$(1.1) \quad \min_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} F^\sigma(\mu), \text{ with } F^\sigma(\mu) := F(\mu) + \sigma \text{KL}(\mu|\pi),$$

over the Wasserstein space $(\mathcal{P}_2(\mathbb{R}^d), \mathcal{W}_2)$, where $F : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ is a function bounded below on $\mathcal{P}_2(\mathbb{R}^d)$, i.e., $\inf_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} F(\mu) > -\infty$, $\pi \in \mathcal{P}_2(\mathbb{R}^d)$ is a reference probability measure, $\sigma > 0$ is a regularization parameter and KL is the KL-divergence (relative entropy). Such optimization problems are motivated by many applications in data science and machine learning, including the task of training two-layer neural networks (NNs) in the mean-field regime [26, 13, 24, 30, 34], generative modeling [18, 2] and reinforcement learning [22, 40].

In this work, we tackle (1.1) from the perspective of discrete-time stepping schemes by proposing the following Jordan–Kinderlehrer–Otto (JKO)-based optimization methods: proximal point, prox-linear and proximal gradient,¹ for which we prove linear convergence to the minimizer of F^σ , without requiring that F is geodesically convex. For $\pi \propto e^{-U}$ with a sufficiently regular potential $U : \mathbb{R}^d \rightarrow \mathbb{R}$ and $F = 0$, given a step-size $\tau > 0$ and starting from $\mu^0 \in \mathcal{P}_2(\mathbb{R}^d)$, the JKO scheme, also called minimizing movement scheme or proximal descent in the Wasserstein space, was originally introduced in [20] and constructs a sequence $(\mu^n)_{n \in \mathbb{N}} \subset \mathcal{P}_2(\mathbb{R}^d)$ by the update rule

$$(1.2) \quad \mu^{n+1} = \operatorname{argmin}_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \left\{ \sigma \text{KL}(\mu|\pi) + \frac{1}{2\tau} \mathcal{W}_2^2(\mu, \mu^n) \right\},$$

2020 *Mathematics Subject Classification.* 46N10, 49Q22, 49K30, 58E30.

Key words and phrases. Entropy regularization, Proximal JKO-based schemes, Optimal transport, Mean-field optimization, logarithmic Sobolev inequality.

¹We maintain the terminology used for analogous methods in finite-dimensional optimization; see e.g. [14, 29]. Indeed, our naming convention is justified since the JKO step (1.2) can be viewed as a proximal operator on the Wasserstein space $(\mathcal{P}_2(\mathbb{R}^d), \mathcal{W}_2)$.

where \mathcal{W}_2 is the L^2 -Wasserstein distance. Note that, by replacing π in (1.2) with $e^{-\sigma^{-1}f-U}$ for some function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, and normalizing appropriately, (1.2) covers optimization problems for functions of the form $F^\sigma(\mu) = \int_{\mathbb{R}^d} f d\mu + \sigma \text{KL}(\mu|e^{-U})$.

Numerical methods for implementing the JKO scheme (1.2) were proposed in [19, 4]. A survey of these methods is also included in [33, Section 4.7]. Moreover, if the initial measure μ^0 and the target measure π are both Gaussian, it is showed in [38, Section 3; Example 5] that the update step in (1.2) can be computed in closed form. Note that, however, these works do not cover the case of non-linear F .

1.1. JKO-based stepping schemes. A natural approach to solving (1.1) for a general F is to start with the proximal point scheme

$$(1.3) \quad \mu^{n+1} = \operatorname{argmin}_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \left\{ F(\mu) + \sigma \text{KL}(\mu|\pi) + \frac{1}{2\tau} \mathcal{W}_2^2(\mu, \mu^n) \right\}.$$

However, this scheme requires one to solve, at each step, a convex but nonlinear minimization problem and hence is mostly of theoretical interest. Its advantage on the theoretical level is that it lends itself to a clean convergence proof with fewest regularity assumptions, which is why we include analysis of this scheme.

A more practical scheme can be created by linearizing F around μ^n and leveraging the fact that the Wasserstein penalty term $\mathcal{W}_2^2(\mu, \mu^n)$ ensures that the linearization is accurate enough, provided there is appropriate regularity of F . Thus, we define the prox-linear scheme

$$(1.4) \quad \mu^{n+1} = \operatorname{argmin}_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \left\{ \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu^n, x)(\mu - \mu^n)(dx) + \sigma \text{KL}(\mu|\pi) + \frac{1}{2\tau} \mathcal{W}_2^2(\mu, \mu^n) \right\}.$$

Since the map $\mu \mapsto \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu^n, x)\mu(dx)$ is linear, by replacing π in (1.2) with appropriately normalized $e^{-\sigma^{-1} \frac{\delta F}{\delta \mu}(\mu^n, \cdot) - U}$, we see that (1.2) covers (1.4) as a special case (cf. the remark below (1.2)). In other words, one could view (1.4) as corresponding to (1.2) with a relative entropy of the form $\text{KL}(\mu|\Phi[\mu^n])$, where $\Phi[\mu^n] \propto e^{-\sigma^{-1} \frac{\delta F}{\delta \mu}(\mu^n, \cdot) - U}$. Thus, (1.4) can be implemented numerically as discussed in [33, Section 4.7]. More recently, [35] proposed an algorithm for solving (1.4) in the case where $F(\mu) = \frac{1}{2} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} W(x, x')\mu(dx)\mu(dx')$, for an interactive potential $W : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ that satisfies $W(x, x') = W(x', x)$, for all $x, x' \in \mathbb{R}^d$. Several numerical experiments were performed but no convergence rates for the algorithm were proved.

Another natural approach is to consider the proximal gradient algorithm, which in our context translates to updating μ^n by a pushforward, which in fact we will show is an optimal transport map for F regular enough and sufficiently small τ , and updating the resulting measure via a JKO step. Thus, the proximal gradient scheme is

$$(1.5) \quad \begin{aligned} \nu^{n+1} &= (I_d - \tau \nabla_{\mu} F(\mu^n)(\cdot))_{\#} \mu^n, \\ \mu^{n+1} &= \operatorname{argmin}_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \left\{ \sigma \text{KL}(\mu|\pi) + \frac{1}{2\tau} \mathcal{W}_2^2(\mu, \nu^{n+1}) \right\}. \end{aligned}$$

A proximal scheme related to (1.5) was recently introduced in [31] for the case where $F = 0$ in (1.1). This method splits $\text{KL}(\mu|\pi)$ into the sum of $\int_{\mathbb{R}^d} U d\mu$ and the entropy $H(\mu)$ (cf. Subsection 2.1) and then implements a gradient descent step on U and a JKO update step for H , see the discussion in Section 1.4 for more details.

It is worth emphasizing that an “explicit” scheme in which both F and $\text{KL}(\cdot|\pi)$ are linearized around μ^n is not expected to converge due to the non-smoothness of the relative

entropy in the Wasserstein space [38, Subsection 3.1.1]. Recently, [39] provided two counterexamples for which updating μ^n by the pushforward

$$\mu^{n+1} = (I_d - \tau \nabla_{\mu} \text{KL}(\mu^n | \pi))_{\#} \mu^n,$$

fails to converge for μ^0 and U appropriately chosen.

Utilizing techniques from optimal transport and the theory of gradient flows on the space of probability measures, we prove that the iterates $(\mu^n)_{n \in \mathbb{N}}$ generated by each of the schemes (1.3), (1.4) and (1.5) converge linearly to the minimizer of F^σ . A notable aspect of our proof is the application of the uniform LSI and a “sandwich” entropy lemma, which makes our work a discrete-time counterpart to [27, 12].

1.2. Connection to the Wasserstein gradient flow. As $\tau \rightarrow 0$, schemes (1.3), (1.4) and (1.5) are expected to recover the Wasserstein gradient flow of F^σ , given by

$$(1.6) \quad \partial_t \mu = \nabla \cdot \left(\left(\nabla \frac{\delta F}{\delta \mu}(\mu, \cdot) + \sigma \nabla U \right) \mu \right) + \sigma \Delta \mu, \quad \mu|_{t=0} := \mu^0 \in \mathcal{P}_2(\mathbb{R}^d).$$

In continuous time, there are two potential approaches to show that (1.6) converges with rate $\mathcal{O}(e^{-\kappa t})$, for $\kappa > 0$, to the minimizer μ_σ^* of F^σ . The appropriate approach depends on F .

On the one hand, assume that F is geodesically convex and U is β -strongly-convex for $\beta > 0$. Then F^σ is $\sigma\beta$ -geodesically convex, which implies that

$$F^\sigma(\mu_\sigma^*) - F^\sigma(\mu_t) \geq \left\langle \nabla \frac{\delta F^\sigma}{\delta \mu}(\mu_t, \cdot), T_{\mu_t}^{\mu_\sigma^*} - I_d \right\rangle_{L_{\mu_t}^2(\mathbb{R}^d)} + \frac{\sigma\beta}{2} \mathcal{W}_2^2(\mu_t, \mu_\sigma^*),$$

where $T_{\mu_t}^{\mu_\sigma^*} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the optimal transport map from μ_t to μ_σ^* , provided that it exists. Furthermore, by [1, Lemma 8.4.7] applied to (1.6), it holds that

$$\frac{1}{2} \frac{d}{dt} \mathcal{W}_2^2(\mu_t, \mu_\sigma^*) = \left\langle \nabla \frac{\delta F^\sigma}{\delta \mu}(\mu_t, \cdot), T_{\mu_t}^{\mu_\sigma^*} - I_d \right\rangle_{L_{\mu_t}^2(\mathbb{R}^d)}.$$

Hence one obtains the following Evolution Variational Inequality (EVI, cf. [1, Theorem 11.1.4])

$$\frac{1}{2} \frac{d}{dt} \mathcal{W}_2^2(\mu_t, \mu_\sigma^*) \leq -(F^\sigma(\mu_t) - F^\sigma(\mu_\sigma^*)) - \frac{\sigma\beta}{2} \mathcal{W}_2^2(\mu_t, \mu_\sigma^*),$$

which implies convergence of (1.6) to μ_σ^* in the Wasserstein distance with rate $\mathcal{O}(e^{-\sigma\beta t})$.

On the other hand, assume that F is flat-convex, which implies that

$$(1.7) \quad F(\mu_\sigma^*) - F(\mu_t) \geq \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu_t, x) (\mu_\sigma^* - \mu_t)(dx),$$

and assume that the proximal measure $\Phi[\mu] \propto e^{-\sigma^{-1} \frac{\delta F}{\delta \mu}(\mu, \cdot) - U}$ satisfies the log-Sobolev inequality (LSI) with a constant $\theta > 0$ for any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$. By (1.7), we obtain

$$(1.8) \quad F^\sigma(\mu_t) - F^\sigma(\mu_\sigma^*) \leq \sigma \text{KL}(\mu_t | \Phi[\mu_t]),$$

which is the right-hand side of the “sandwich” entropy lemma (Lemma 2.7). Then via the arguments in [27, 12] using the LSI and (1.8), we have

$$\begin{aligned} \frac{d}{dt} (F^\sigma(\mu_t) - F^\sigma(\mu_\sigma^*)) &= -\sigma^2 I(\mu_t | \Phi[\mu_t]) \leq -2\theta\sigma^2 \text{KL}(\mu_t | \Phi[\mu_t]) \\ &\leq -2\theta\sigma (F^\sigma(\mu_t) - F^\sigma(\mu_\sigma^*)). \end{aligned}$$

Hence, convergence of (1.6) to μ_σ^* with rate $\mathcal{O}(e^{-2\theta\sigma t})$ is obtained from Gronwall’s lemma. In this setting, the same rate of convergence in the Wasserstein distance then follows immediately from Lemma 2.7 and Talagrand’s inequality since $\mu_\sigma^* = \Phi[\mu_\sigma^*]$.

We stress that the proof strategy via EVI fails if F is not geodesically convex, and that there are examples of applications where the assumption of geodesic convexity is not satisfied, while flat convexity holds (cf. Example 3.3). Motivated by this, in the present paper we adapt the proof via LSI and the “sandwich” entropy lemma to discrete-time stepping schemes, and prove linear convergence of (1.3), (1.4) and (1.5) to μ_σ^* . In particular, our proof only requires the notion of flat convexity of F instead of more restrictive geodesic convexity.

A related line of research focuses on establishing LSIs for particle approximations of the mean-field Langevin dynamics (1.6). This has received significant attention lately with positive results [11, 37, 25].

1.3. Our contribution. We propose JKO-based methods for solving the mean-field optimization problem (1.1). Our contribution can be summarized as follows:

- In Theorem 5.1, 6.1, 7.2, we prove existence and uniqueness of the minimizer for each scheme (1.3), (1.4) and (1.5), respectively.
- In Proposition 5.3, 6.3, 7.4, we prove that along the iterates generated by these schemes the relative entropy $\text{KL}(\cdot|\pi)$ admits a unique Wasserstein subgradient, and hence, in Lemma 5.4, 6.4, 7.5, we show that the iterates satisfy first-order optimality conditions.
- Finally, our main contributions are Theorem 3.1 and Corollary 3.2 where we prove linear convergence of (1.3), (1.4) and (1.5) to the minimizer μ_σ^* of F^σ , with respect to $F^\sigma(\cdot) - F^\sigma(\mu_\sigma^*)$, $\text{KL}(\cdot|\mu_\sigma^*)$ and $\mathcal{W}_2^2(\cdot, \mu_\sigma^*)$. In particular, we show that for each of the schemes there exists $\kappa > 1$ such that

$$0 \leq F^\sigma(\mu^n) - F^\sigma(\mu_\sigma^*) \leq \kappa^{-n} (F^\sigma(\mu^0) - F^\sigma(\mu_\sigma^*)),$$

for all $n \in \mathbb{N}$.

1.4. Related works. As discussed in the first part of the introduction, [31] considered the case where $F = 0$, $\sigma = 1$ and $\pi \propto e^{-U}$ in (1.1), with U strongly convex, L_U -smooth for some $L_U > 0$, and the proximal scheme

$$(1.9) \quad \begin{aligned} \nu^{n+1} &= (I_d - \tau \nabla U)_{\#} \mu^n, \\ \mu^{n+1} &= \operatorname{argmin}_{\mu \in \mathcal{P}_2^\lambda(\mathbb{R}^d)} \left\{ H(\mu) + \frac{1}{2\tau} \mathcal{W}_2^2(\mu, \nu^{n+1}) \right\}, \end{aligned}$$

with step-size $\tau > 0$ and starting from $\mu^0 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, where H is the entropy (cf. (2.1)). Given that $\tau < L_U^{-1}$, it is proved in [31, Corollary 11] that the iterates $(\mu^n)_{n \in \mathbb{N}}$ generated by (1.9) converge with linear rate with respect to $\mu \mapsto \mathcal{W}_2^2(\mu, \pi)$ by establishing a discrete-time variant of the EVI (see [31, Proposition 8]). One of the key ingredients in [31] for proving the EVI is the geodesic convexity of $\int_{\mathbb{R}^d} U d\mu$, which follows from assuming strong convexity of U .

As we have already mentioned, for problem (1.1) the approach via EVI fails since F is not necessarily geodesically convex. Hence, we prove convergence of (1.3), (1.4) and (1.5) with linear rate via the LSI and the entropy “sandwich” lemma. We work under the assumptions that F is flat-convex and U is strongly-convex. Moreover, since the JKO step in (1.5) uses the relative entropy $\text{KL}(\cdot|\pi)$ instead of the entropy H , L_U -smoothness of U is not needed.

The setting of [31] was extended in [23] by assuming that U is non-convex but expressed as $U := U_1 - U_2$, where $U_1, U_2 : \mathbb{R}^d \rightarrow \mathbb{R}$ are convex, have quadratic growth and U_2 is L_{U_2} -smooth, for some $L_{U_2} > 0$. In this setup, the proximal scheme considered in [23] is

$$(1.10) \quad \begin{aligned} \nu^{n+1} &= (I_d + \tau \nabla U_2)_{\#} \mu^n, \\ \mu^{n+1} &= \operatorname{argmin}_{\mu \in \mathcal{P}_2^\lambda(\mathbb{R}^d)} \left\{ \int_{\mathbb{R}^d} U_1(x) dx + H(\mu) + \frac{1}{2\tau} \mathcal{W}_2^2(\mu, \nu^{n+1}) \right\}, \end{aligned}$$

with step-size $\tau > 0$ and starting from $\mu^0 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$. Under the assumption that $\pi \propto e^{-(U_1 - U_2)}$ satisfies LSI, it is proved in [23, Theorem 4.5 (ii)] that the iterates $(\mu^n)_{n \in \mathbb{N}}$ generated by (1.10) converge with linear rate with respect to $\mu \mapsto \operatorname{KL}(\mu | \pi)$ and $\mu \mapsto \mathcal{W}_2^2(\mu, \pi)$. However, to apply the results from [23], one would need to verify that $\pi \propto e^{-(U_1 - U_2)}$ indeed satisfies LSI, which is not immediate, and may require additional conditions on U_1 and U_2 . Although in our setting the potential $\sigma^{-1} \frac{\delta F}{\delta \mu}(\mu, \cdot) + U$ is also non-convex, we rigorously show that the proximal measure $\Phi[\mu] \propto e^{-\sigma^{-1} \frac{\delta F}{\delta \mu}(\mu, \cdot) - U}$ satisfies LSI for any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ due to the Holley–Stroock criterion.

Moreover, the existence and uniqueness of the minimizer for the JKO step in (1.9) and (1.10) are just postulated as assumptions in both [31, Assumption B2] and [23, Assumption 1 (iv)], respectively. In contrast, we provide a fully rigorous proof of the existence and uniqueness of the minimizer for each of our schemes (1.3), (1.4) and (1.5).

We also cover the issue of Wasserstein sub-differentiability of the relative entropy $\operatorname{KL}(\cdot | \pi)$ at every step of each scheme (1.3), (1.4) and (1.5). Since the relative entropy is not Wasserstein differentiable (but only Wasserstein sub-differentiable), one needs to prove that along each scheme the iterates belong to a certain class of Sobolev regularity, and hence the relative entropy $\operatorname{KL}(\cdot | \pi)$ has a unique Wasserstein sub-gradient at μ given by $\nabla \log \frac{d\mu}{d\pi}$. Moreover, within this class we are guaranteed that the relative Fisher information is finite. These two points, namely existence of unique minimizers for scheme steps and Wasserstein differentiability are in fact technically the most challenging steps of all our convergence proofs.

2. PRELIMINARIES AND ASSUMPTIONS

In this section, we introduce the necessary notations and assumptions used throughout the paper, and recall some definitions and results.

2.1. Notation. By $\mathcal{P}_2(\mathbb{R}^d)$ we denote the space of probability measures with finite second moment. We equip $\mathcal{P}_2(\mathbb{R}^d)$ with the 2-Wasserstein distance \mathcal{W}_2 , and as it is common in the literature we refer to the metric space $(\mathcal{P}_2(\mathbb{R}^d), \mathcal{W}_2)$ as the Wasserstein space. Let $\mathcal{B}(\mathbb{R}^d)$ denote the Borel σ -algebra over \mathbb{R}^d . For any measures μ, ν on $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$, we write $\mu \ll \nu$ if μ is absolutely continuous with respect to ν . The set of absolutely continuous measures in $\mathcal{P}_2(\mathbb{R}^d)$ with respect to ν is denoted by $\mathcal{P}_2^\nu(\mathbb{R}^d) := \{\mu \in \mathcal{P}_2(\mathbb{R}^d) : \mu \ll \nu\}$. We denote by λ the Lebesgue measure on \mathbb{R}^d . For any measures $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$, the map $T_\mu^\nu : \mathbb{R}^d \rightarrow \mathbb{R}^d$ denotes the optimal transport map from μ to ν . Let $p \in \{1, 2\}$. For any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, let $(L_\mu^p(\mathbb{R}^d), \|\cdot\|_{L_\mu^p(\mathbb{R}^d)})$ be the space of $\mathcal{B}(\mathbb{R}^d)$ -measurable functions $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that $\|f\|_{L_\mu^p(\mathbb{R}^d)} := \left(\int_{\mathbb{R}^d} |f(x)|^p \mu(dx) \right)^{\frac{1}{p}} < \infty$. Note that the identity map $I_d : \mathbb{R}^d \rightarrow \mathbb{R}^d$, given by $I(x) = x$, for all $x \in \mathbb{R}^d$, is an element of $L_\mu^2(\mathbb{R}^d)$. For any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, we denote by $\langle \cdot, \cdot \rangle_{L_\mu^2(\mathbb{R}^d)}$ the inner product on the space $L_\mu^2(\mathbb{R}^d)$. Let $W_\mu^{1,p}(\mathbb{R}^d)$ be the weighted Sobolev space of $\mathcal{B}(\mathbb{R}^d)$ -measurable functions $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$

such that $f \in L^p_\mu(\mathbb{R}^d)$ and $\nabla f \in L^p_\mu(\mathbb{R}^d)$. Let $W_{\lambda, \text{loc}}^{1,1}(\mathbb{R}^d)$ be the Sobolev space of $\mathcal{B}(\mathbb{R}^d)$ -measurable functions $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that $f \in L^1_{\lambda, \text{loc}}(\mathbb{R}^d)$ and $\nabla f \in L^1_{\lambda, \text{loc}}(\mathbb{R}^d)$. For any $f, g : \mathbb{R}^d \rightarrow \mathbb{R}^d$, the composition $f \circ g : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is denoted by $f(g)$.

For the Lebesgue measure λ on \mathbb{R}^d , the entropy $H : \mathcal{P}_2(\mathbb{R}^d) \rightarrow [0, \infty]$ is given for any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ by

$$(2.1) \quad H(\mu) := \begin{cases} \int_{\mathbb{R}^d} \log \frac{d\mu}{d\lambda}(x) \mu(dx), & \mu \in \mathcal{P}_2^\lambda(\mathbb{R}^d), \\ +\infty, & \text{else.} \end{cases}$$

For $\pi \in \mathcal{P}_2(\mathbb{R}^d)$, the relative entropy $\text{KL}(\cdot|\pi) : \mathcal{P}_2(\mathbb{R}^d) \rightarrow [0, \infty]$ with respect to π is given for any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ by

$$\text{KL}(\mu|\pi) := \begin{cases} \int_{\mathbb{R}^d} \log \frac{d\mu}{d\pi}(x) \mu(dx), & \mu \in \mathcal{P}_2^\pi(\mathbb{R}^d), \\ +\infty, & \text{else,} \end{cases}$$

and the relative Fisher information $I(\cdot|\pi) : \mathcal{P}_2(\mathbb{R}^d) \rightarrow [0, \infty]$ with respect to π is given for any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ by

$$I(\mu|\pi) := \begin{cases} \int_{\mathbb{R}^d} \left| \nabla \log \frac{d\mu}{d\pi}(x) \right|^2 \mu(dx), & \mu \in \mathcal{P}_2^\pi(\mathbb{R}^d) \text{ and } \sqrt{\frac{d\mu}{d\pi}} \in W_\pi^{1,2}(\mathbb{R}^d), \\ +\infty, & \text{else.} \end{cases}$$

Assumption 2.1 (Flat-convexity of F). Assume that $F \in \mathcal{C}^1$ (cf. Definition B.1) is convex on $\mathcal{P}_2(\mathbb{R}^d)$, i.e., for any $\mu', \mu \in \mathcal{P}_2(\mathbb{R}^d)$, it holds

$$F(\mu') - F(\mu) \geq \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu, x) (\mu' - \mu)(dx).$$

Assumption 2.2 (Lipschitzness and boundedness of the flat derivative). Assume that $F \in \mathcal{C}^1$ and there exist $C_F, L_F > 0$ such that for all $\mu, \mu' \in \mathcal{P}_2(\mathbb{R}^d)$ and all $x, x' \in \mathbb{R}^d$, we have

$$(2.2) \quad \left| \frac{\delta F}{\delta \mu}(\mu', x') - \frac{\delta F}{\delta \mu}(\mu, x) \right| \leq L_F (|x' - x| + \mathcal{W}_2(\mu', \mu)),$$

$$(2.3) \quad \left| \frac{\delta F}{\delta \mu}(\mu, x) \right| \leq C_F.$$

Assumption 2.3. Assume that $\pi(dx) \propto e^{-U(x)} dx$ for a continuously differentiable function $U : \mathbb{R}^d \rightarrow \mathbb{R}$ such that:

(i) U is bounded below and has at least quadratic growth, i.e.,

$$(2.4) \quad \text{ess inf}_{x \in \mathbb{R}^d} U(x) > -\infty \text{ and } \liminf_{x \rightarrow \infty} \frac{U(x)}{|x|^2} > 0,$$

(ii) U is α_U -strongly convex, i.e., there exists $\alpha_U > 0$ such that for all $x, y \in \mathbb{R}^d$, it holds

$$(2.5) \quad \alpha_U |x - y|^2 \leq (x - y) \cdot (\nabla U(x) - \nabla U(y)),$$

Assumption 2.1, 2.2 and 2.3 are standard in the mean-field optimization literature; see e.g. [17, 27, 12, 9]. In particular, the last two allow us to establish existence and uniqueness of the minimizer of (1.1).

Proposition 2.4 ([9, Proposition 1]). *Let Assumption 2.2 and (2.4) in Assumption 2.3 hold. Then F^σ admits a unique minimizer $\mu_\sigma^* \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$ given by*

$$\mu_\sigma^*(dx) = \frac{1}{Z(\mu_\sigma^*)} e^{-\frac{1}{\sigma} \frac{\delta F}{\delta \mu}(\mu_\sigma^*, x) - U(x)} dx,$$

where $Z(\mu_\sigma^*)$ is a normalization constant.

Next, we recall the definition of the so-called proximal Gibbs measure which is a crucial ingredient in proving convergence via LSI.

Definition 2.5 (Proximal Gibbs measure; [27, 12]). For any $\mu \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, define the operator $\Phi : \mathcal{P}_2^\lambda(\mathbb{R}^d) \rightarrow \mathcal{P}_2^\lambda(\mathbb{R}^d)$ by

$$(2.6) \quad \Phi[\mu](dx) := \frac{1}{Z(\mu)} e^{-\frac{1}{\sigma} \frac{\delta F}{\delta \mu}(\mu, x) - U(x)} dx,$$

where, for each $\mu \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, $Z(\mu)$ is a normalization constant. We call $\Phi[\mu]$ the proximal Gibbs measure.

It follows via (2.5) in Assumption 2.3 and the Bakry–Émery criterion [3], that π satisfies the LSI with constant $\alpha_U > 0$. This fact together with (2.3) in Assumption 2.2 and the Holley–Stroock criterion [16] imply that $\Phi[\mu]$ satisfies the LSI with constant $\alpha_U e^{-\frac{4C_F}{\sigma}}$, i.e., for any $\mu \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, it holds

$$(2.7) \quad \text{KL}(\mu | \Phi[\mu]) \leq \frac{e^{\frac{4C_F}{\sigma}}}{2\alpha_U} I(\mu | \Phi[\mu]).$$

Furthermore, according to [28], since $\Phi[\mu]$ satisfies (2.7), it also satisfies Talagrand’s inequality, i.e., for any $\mu \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, it holds

$$(2.8) \quad \mathcal{W}_2^2(\mu, \Phi[\mu]) \leq \frac{2e^{\frac{4C_F}{\sigma}}}{\alpha_U} \text{KL}(\mu | \Phi[\mu]).$$

We also highlight that μ_σ^* satisfies both (2.7) and (2.8) since $\mu_\sigma^* = \Phi[\mu_\sigma^*]$.

Remark 2.6. Condition (2.3) in Assumption 2.2 that is used for the Holley–Stroock criterion could be relaxed into Lipschitz continuity of the map $x \mapsto \frac{\delta F}{\delta \mu}(\mu, x)$ uniformly over μ . Under this assumption, according to [8, Theorem 2.7 (2)], we obtain an upper bound on the log-Sobolev constant of $\Phi[\mu]$. Thus, our results will still hold if we replace $\alpha_U e^{-\frac{4C_F}{\sigma}}$ by that upper bound.

Using Assumption 2.1, [27, 12] proved the following entropy “sandwich” lemma, which provides bounds for the distance between F^σ and the minimum value $F^\sigma(\mu_\sigma^*)$.

Lemma 2.7 ([27, 12]). *Let Assumption 2.1, 2.2 and (2.4) in Assumption 2.3 hold. Let $\mu_\sigma^* \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$ be the unique minimizer of F^σ . Then, for any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, we have*

$$\sigma \text{KL}(\mu | \mu_\sigma^*) \leq F^\sigma(\mu) - F^\sigma(\mu_\sigma^*) \leq \sigma \text{KL}(\mu | \Phi[\mu]).$$

It is worth stressing that for proving linear convergence of (1.3), (1.4) and (1.5), it suffices to show that for each scheme, there exists $\kappa > 1$ such that

$$F^\sigma(\mu^n) - F^\sigma(\mu_\sigma^*) \leq \kappa^{-n} (F^\sigma(\mu^0) - F^\sigma(\mu_\sigma^*)),$$

for all $n \in \mathbb{N}$, where $(\mu^n)_n$ are the iterates generated by (1.3), (1.4) and (1.5), respectively. Then convergence with respect to $\mu \mapsto \text{KL}(\mu | \mu_\sigma^*)$ and $\mu \mapsto \mathcal{W}_2^2(\mu, \mu_\sigma^*)$ will follow immediately from Lemma 2.7 and Talagrand’s inequality (2.8). The dependence of other constants on κ will be made explicit in the statement of the convergence result.

The following standard assumptions are concerned with the Wasserstein regularity of F , in particular the existence of its Wasserstein gradient understood as the Euclidean gradient of the flat derivative, and the Lipschitz continuity of the Wasserstein gradient.

Assumption 2.8 (Wasserstein differentiability of F). Assume that, for any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, the function $\mathbb{R}^d \ni x \mapsto \frac{\delta F}{\delta \mu}(\mu, x) \in \mathbb{R}$ is differentiable,

- (i) there exists $C'_F > 0$ such that $\left| \nabla \frac{\delta F}{\delta \mu}(\mu, x) \right| \leq C'_F$, for all $(\mu, x) \in \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d$, and
- (ii) the derivative $\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d \ni (\mu, x) \mapsto \nabla \frac{\delta F}{\delta \mu}(\mu, x) \in \mathbb{R}^d$ is jointly continuous in (μ, x) .

Under Assumption 2.8, by [7, Proposition 5.48; Theorem 5.64], it follows that $F : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ is Wasserstein differentiable at $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ (cf. Definition B.4) and $\nabla_\mu F(\mu)(\cdot) = \nabla \frac{\delta F}{\delta \mu}(\mu, \cdot)$.

Assumption 2.9 (Lipschitzness of the Wasserstein gradient). Assume there exists $L'_F > 0$ such that for all $\mu, \mu' \in \mathcal{P}_2(\mathbb{R}^d)$ and all $x, x' \in \mathbb{R}^d$, we have

$$|\nabla_\mu F(\mu')(x') - \nabla_\mu F(\mu)(x)| \leq L'_F (|x' - x| + \mathcal{W}_2(\mu', \mu)).$$

3. MAIN RESULTS

In this section, we present the main convergence results. There are three groups of auxiliary results on which the convergence proofs are built, each corresponding to one of the schemes (1.3), (1.4) and (1.5), but these are deferred to Appendix 5, 6 and 7, respectively. Before we state the main results, we first outline the general proof steps applicable to all schemes:

- (1) Modifying the argument in the proof of [20, Proposition 4.1], we prove that given $\mu^n \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, each scheme (1.3), (1.4) and (1.5) admits a unique minimizer $\mu^{n+1} \in \mathcal{P}_2^\pi(\mathbb{R}^d)$.
- (2) Then, leveraging results that connect the metric slope of the relative entropy and the relative Fisher information, we prove that, for each scheme (1.3), (1.4) and (1.5),

$$\mu^{n+1} \in \mathfrak{C} := \left\{ m \in \mathcal{P}_2^\pi(\mathbb{R}^d) : \frac{dm}{d\pi} \in W_{\lambda, \text{loc}}^{1,1}(\mathbb{R}^d), \sqrt{\frac{dm}{d\pi}} \in W_\pi^{1,2}(\mathbb{R}^d) \right\}.$$

Hence, by Theorem A.5, we conclude that the relative entropy $\text{KL}(\cdot|\pi)$ has a unique Wasserstein subgradient at μ^{n+1} , and it is given by $\nabla \log \frac{d\mu^{n+1}}{d\pi}$.

- (3) As a consequence of the previous result, we prove first-order optimality conditions for each scheme (1.3), (1.4) and (1.5). The optimality conditions enable us to connect the Fisher information relative to the proximal measures Φ given by (2.6) with the Wasserstein distance $\mathcal{W}_2^2(\mu^{n+1}, \mu^n)$.
- (4) Finally, for (1.3), the convergence proof is concluded via the LSI and Lemma 2.7. For (1.4) and (1.5), the proofs are concluded in the same way as for (1.3) but they also require smoothness of F relative to $\mathcal{W}_2^2(\mu^{n+1}, \mu^n)$, and convexity along (generalized) geodesics of $\text{KL}(\cdot|\pi)$.

Theorem 3.1 (Linear convergence of the schemes (1.3), (1.4), 1.5). *Let Assumption 2.1, 2.2, 2.3, 2.8 hold. Let $\mu^0 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$. Then*

- (i) *for the iterates $(\mu^n)_{n \geq 1}$ of (1.3), we have*

$$F^\sigma(\mu^n) - F^\sigma(\mu_\sigma^*) \leq \left(1 + e^{-\frac{4C_F}{\sigma}} \tau \sigma \alpha_U \right)^{-n} (F^\sigma(\mu^0) - F^\sigma(\mu_\sigma^*)), \text{ for all } n \in \mathbb{N}.$$

- (ii) *Let Assumption 2.9 hold. If $\tau < \frac{2}{L'_F}$, then for the iterates $(\mu^n)_{n \geq 1}$ of (1.4), we have*

$$F^\sigma(\mu^n) - F^\sigma(\mu_\sigma^*) \leq \left(1 + \frac{\tau \sigma \alpha_U (1 - 2\tau L'_F)}{1 + (\tau L'_F)^2} e^{-\frac{4C_F}{\sigma}} \right)^{-n} (F^\sigma(\mu^0) - F^\sigma(\mu_\sigma^*)), \text{ for all } n \in \mathbb{N}.$$

(iii) Let Assumption 2.9 hold. If $\tau < \frac{1}{L'_F}$, then for the iterates $(\mu^n)_{n \geq 1}$ of (1.5), we have

$$F^\sigma(\mu^n) - F^\sigma(\mu_\sigma^*) \leq \left(1 + \frac{\tau\sigma\alpha_U(1 - \tau L'_F)}{1 + 4(\tau L'_F)^2} e^{-\frac{4C_F}{\sigma}}\right)^{-n} (F^\sigma(\mu^0) - F^\sigma(\mu_\sigma^*)), \text{ for all } n \in \mathbb{N}.$$

As mentioned in Section 2, once we have Theorem 3.1, the convergence with respect to $\mu \mapsto \text{KL}(\mu|\mu_\sigma^*)$ and $\mu \mapsto \mathcal{W}_2^2(\mu, \mu_\sigma^*)$ follows from Lemma 2.7 and the Talagrand inequality (2.8).

Corollary 3.2 (Linear convergence in KL and \mathcal{W}_2^2). *For $\kappa > 1$ corresponding to each rate in (i), (ii) and (iii) in Theorem 3.1, we obtain*

$$\begin{aligned} \text{KL}(\mu^n|\mu_\sigma^*) &\leq \frac{1}{\sigma} \kappa^{-n} (F^\sigma(\mu^0) - F^\sigma(\mu_\sigma^*)) \text{ and} \\ \mathcal{W}_2^2(\mu^n, \mu_\sigma^*) &\leq \frac{2e^{\frac{4C_F}{\sigma}}}{\alpha_U \sigma} \kappa^{-n} (F^\sigma(\mu^0) - F^\sigma(\mu_\sigma^*)), \text{ for all } n \in \mathbb{N}. \end{aligned}$$

We finish this section by considering the example of an L^2 -loss function for a two-layer mean-field NN and showing that, under regularity conditions on the activation function, the loss function satisfies Assumption 2.1, 2.2, 2.8, 2.9.

Example 3.3 (Two-layer mean-field neural network; [17]). Let ν be a compactly supported measure representing the training data $(y, z) \in \mathbb{R} \times \mathbb{R}^{d-1}$, let $(w, b) \in \mathbb{R}^{d-1} \times \mathbb{R}$ be the parameters of the neural network and let $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ be a bounded, continuous, non-constant activation function.

For $x := (w, b) \in \mathbb{R}^d$ and $z \in \mathbb{R}^{d-1}$, define the function $\hat{\varphi}(x, z) := \ell(b)\varphi(\langle w, z \rangle)$, where $\ell : \mathbb{R} \rightarrow [-K, K]$ is a clipping function with clipping threshold $K > 0$. The training of the two-layer neural network aims to find the optimal set of parameters $\{x_i\}_{i=1}^N$ which minimize the non-convex L^2 -loss function

$$(3.1) \quad F_N(x_1, \dots, x_N) := \int_{\mathbb{R}^d} \left| y - \frac{1}{N} \sum_{i=1}^N \hat{\varphi}(x_i, z) \right|^2 \nu(dy, dz).$$

Instead of the non-convex minimization problem (3.1), we consider the mean-field optimization problem

$$\min_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} F(\mu), \quad \text{with } F(\mu) := \int_{\mathbb{R}^d} \left| y - \mathbb{E}^{X \sim \mu}[\hat{\varphi}(X, z)] \right|^2 \nu(dy, dz).$$

Observe that by linearity of the expectation in μ and convexity of $|\cdot|^2$, the function F satisfies the flat-convexity condition $F((1 - \varepsilon)\mu + \varepsilon\mu') \leq (1 - \varepsilon)F(\mu) + \varepsilon F(\mu')$, for any $\mu, \mu' \in \mathcal{P}_2(\mathbb{R}^d)$ and any $\varepsilon \in [0, 1]$. Hence, by [17, Lemma 4.1], F satisfies Assumption 2.1.

We stress that F is not geodesically convex in the Wasserstein space. If it were, then this would require F_N to be convex, which is clearly not the case (see also [10, Remark 3.4] for more details).

Assume that $\nabla_x \hat{\varphi}, \nabla_x^2 \hat{\varphi}, \nabla_b \ell, \nabla_b^2 \ell$ are bounded and continuous. Then

$$\begin{aligned} \frac{\delta F}{\delta \mu}(\mu, x) &= - \int_{\mathbb{R}^d} (y - \mathbb{E}^{X \sim \mu}[\hat{\varphi}(X, z)]) \hat{\varphi}(x, z) \nu(dy, dz), \\ \nabla_\mu F(\mu)(x) &= - \int_{\mathbb{R}^d} (y - \mathbb{E}^{X \sim \mu}[\hat{\varphi}(X, z)]) \nabla_x \hat{\varphi}(x, z) \nu(dy, dz), \end{aligned}$$

and a straightforward calculation shows that F satisfies Assumption 2.2, 2.8, 2.9. In view of Remark (2.6), note that the clipping function ℓ is not needed if Lipschitz continuity of the map $x \mapsto \frac{\delta F}{\delta \mu}(\mu, x)$ uniformly over μ is assumed instead of condition (2.3) in Assumption 2.2.

Before we give the proof of Theorem 3.1, for convenience, we give a short calculation that we will repeatedly use in the proof. For F satisfying Assumption 2.8, $\pi \propto e^{-U}$, any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ and any $\mu' \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, it holds

$$(3.2) \quad \begin{aligned} \nabla_\mu F(\mu) + \sigma \nabla \log \frac{d\mu'}{d\pi} &= \nabla \frac{\delta F}{\delta \mu}(\mu, \cdot) + \sigma \nabla \log \frac{d\mu'}{d\pi} \\ &= -\sigma \nabla \log e^{-\frac{1}{\sigma} \frac{\delta F}{\delta \mu}(\mu, \cdot)} + \sigma \nabla \log \frac{d\mu'}{d\pi} = \sigma \nabla \log \frac{d\mu'}{d\Phi[\mu]}, \end{aligned}$$

where the last equality follows from (2.6).

4. PROOF OF THEOREM 3.1

4.1. Proof of (i). First, note that by Proposition 5.3, the iterates $(\mu^n)_{n \in \mathbb{N}} \subset \mathfrak{C}$. Since $\mu^{n+1} \in \mathfrak{C}$ is a minimizer of (1.3), it follows that

$$(4.1) \quad F^\sigma(\mu^{n+1}) + \frac{1}{2\tau} \mathcal{W}_2^2(\mu^{n+1}, \mu^n) \leq F^\sigma(\mu^n) + \frac{1}{2\tau} \mathcal{W}_2^2(\mu^n, \mu^n) = F^\sigma(\mu^n).$$

By Lemma 5.4, the optimality condition for (1.3) reads

$$\nabla_\mu F(\mu^{n+1})(x) + \sigma \nabla \log \frac{d\mu^{n+1}}{d\pi}(x) = \frac{1}{\tau} \left(T_{\mu^{n+1}}^{\mu^n}(x) - x \right), \quad \text{for } \mu^{n+1}\text{-a.e. } x.$$

By (3.2) with $\mu' = \mu = \mu^{n+1}$, we obtain

$$(4.2) \quad \nabla \log \frac{d\mu^{n+1}}{d\Phi[\mu^{n+1}]}(x) = \frac{1}{\tau\sigma} \left(T_{\mu^{n+1}}^{\mu^n}(x) - x \right), \quad \text{for } \mu^{n+1}\text{-a.e. } x.$$

Squaring both sides of (4.2) and integrating with respect to μ^{n+1} gives

$$(4.3) \quad I \left(\mu^{n+1} \middle| \Phi[\mu^{n+1}] \right) = \frac{1}{\tau^2 \sigma^2} \mathcal{W}_2^2(\mu^{n+1}, \mu^n).$$

Using (4.3) in (4.1) gives

$$\begin{aligned} F^\sigma(\mu^{n+1}) &\leq F^\sigma(\mu^n) - \frac{1}{2\tau} \mathcal{W}_2^2(\mu^{n+1}, \mu^n) = F^\sigma(\mu^n) - \frac{\tau\sigma^2}{2} I \left(\mu^{n+1} \middle| \Phi[\mu^{n+1}] \right) \\ &\leq F^\sigma(\mu^n) - \tau\sigma^2 \frac{\alpha_U}{e^{\frac{4C_F}{\sigma}}} \text{KL} \left(\mu^{n+1} \middle| \Phi[\mu^{n+1}] \right) \leq F^\sigma(\mu^n) - \tau\sigma \frac{\alpha_U}{e^{\frac{4C_F}{\sigma}}} (F^\sigma(\mu^{n+1}) - F^\sigma(\mu_\sigma^*)), \end{aligned}$$

where the second and third inequalities follow from (2.7) and Lemma 2.7, respectively. Let $\kappa := 1 + \frac{\tau\sigma\alpha_U}{e^{\frac{4C_F}{\sigma}}} > 1$. Rearranging the inequality above gives

$$F^\sigma(\mu^{n+1}) - F^\sigma(\mu_\sigma^*) \leq \kappa^{-1} (F^\sigma(\mu^n) - F^\sigma(\mu_\sigma^*)).$$

The convergence estimate follows by iterating over $n \in \mathbb{N}$.

4.2. Proof of (ii). First, note that by Proposition 6.3, the iterates $(\mu^n)_{n \in \mathbb{N}} \subset \mathfrak{C}$. Combining Lemma 6.5 with Lemma 6.6 gives

$$(4.4) \quad \begin{aligned} F^\sigma(\mu^{n+1}) - F^\sigma(\mu^n) &- \left\langle \nabla_\mu F(\mu^n) + \sigma \nabla \log \frac{d\mu^{n+1}}{d\pi} \left(T_{\mu^n}^{\mu^{n+1}} \right), T_{\mu^n}^{\mu^{n+1}} - Id \right\rangle_{L^2_{\mu^n}(\mathbb{R}^d)} \\ &\leq L'_F \mathcal{W}_2^2(\mu^{n+1}, \mu^n). \end{aligned}$$

Using (3.2) with $\mu = \mu^n$ and $\mu' = \mu^{n+1}$, observe that for μ^n -a.e. x , we have

$$\begin{aligned}
(4.5) \quad & \nabla_{\mu} F(\mu^n)(x) + \sigma \nabla \log \frac{d\mu^{n+1}}{d\pi} \left(T_{\mu^n}^{\mu^{n+1}}(x) \right) \\
& = \nabla_{\mu} F(\mu^n)(x) - \nabla_{\mu} F(\mu^n) \left(T_{\mu^n}^{\mu^{n+1}}(x) \right) \\
& \quad + \nabla_{\mu} F(\mu^n) \left(T_{\mu^n}^{\mu^{n+1}}(x) \right) + \sigma \nabla \log \frac{d\mu^{n+1}}{d\pi} \left(T_{\mu^n}^{\mu^{n+1}}(x) \right) \\
& = \nabla_{\mu} F(\mu^n)(x) - \nabla_{\mu} F(\mu^n) \left(T_{\mu^n}^{\mu^{n+1}}(x) \right) + \sigma \nabla \log \frac{d\mu^{n+1}}{d\Phi[\mu^n]} \left(T_{\mu^n}^{\mu^{n+1}}(x) \right).
\end{aligned}$$

Hence, using (4.5) in (4.4) gives

$$\begin{aligned}
(4.6) \quad & F^{\sigma}(\mu^{n+1}) - F^{\sigma}(\mu^n) - \left\langle \sigma \nabla \log \frac{d\mu^{n+1}}{d\Phi[\mu^n]} \left(T_{\mu^n}^{\mu^{n+1}}(x) \right), T_{\mu^n}^{\mu^{n+1}} - I_d \right\rangle_{L_{\mu^n}^2(\mathbb{R}^d)} \\
& \leq L'_F \mathcal{W}_2^2(\mu^{n+1}, \mu^n) + \left\langle \nabla_{\mu} F(\mu^n)(x) - \nabla_{\mu} F(\mu^n) \left(T_{\mu^n}^{\mu^{n+1}}(x) \right), T_{\mu^n}^{\mu^{n+1}} - I_d \right\rangle_{L_{\mu^n}^2(\mathbb{R}^d)} \\
& \leq L'_F \mathcal{W}_2^2(\mu^{n+1}, \mu^n) + \left\| \nabla_{\mu} F(\mu^n) - \nabla_{\mu} F(\mu^n) \left(T_{\mu^n}^{\mu^{n+1}} \right) \right\|_{L_{\mu^n}^2(\mathbb{R}^d)} \left\| T_{\mu^n}^{\mu^{n+1}} - I_d \right\|_{L_{\mu^n}^2(\mathbb{R}^d)} \\
& \leq L'_F \mathcal{W}_2^2(\mu^{n+1}, \mu^n) + L'_F \left\| T_{\mu^n}^{\mu^{n+1}} - I_d \right\|_{L_{\mu^n}^2(\mathbb{R}^d)}^2 = 2L'_F \mathcal{W}_2^2(\mu^{n+1}, \mu^n),
\end{aligned}$$

where the second, third and last inequality follows from the Cauchy-Schwarz inequality, Assumption 2.9 and Corollary A.3, respectively.

By Lemma 6.4, the optimality condition for (1.4) reads

$$\nabla_{\mu} F(\mu^n)(x) + \sigma \nabla \log \frac{d\mu^{n+1}}{d\pi}(x) = \frac{1}{\tau} \left(T_{\mu^{n+1}}^{\mu^n}(x) - x \right), \quad \text{for } \mu^{n+1}\text{-a.e. } x.$$

By Definition 2.5 for $\mu = \mu^n$ and $\mu' = \mu^{n+1}$, we obtain

$$(4.7) \quad \nabla \log \frac{d\mu^{n+1}}{d\Phi[\mu^n]}(x) = \frac{1}{\tau\sigma} \left(T_{\mu^{n+1}}^{\mu^n}(x) - x \right), \quad \text{for } \mu^{n+1}\text{-a.e. } x.$$

Hence, using the fact that $T_{\mu^n}^{\mu^{n+1}} \circ T_{\mu^{n+1}}^{\mu^n} = I_d$, μ^{n+1} -a.e., it follows that

$$(4.8) \quad \nabla \log \frac{d\mu^{n+1}}{d\Phi[\mu^n]} \left(T_{\mu^n}^{\mu^{n+1}}(x) \right) = \frac{1}{\tau\sigma} \left(x - T_{\mu^n}^{\mu^{n+1}}(x) \right), \quad \text{for } \mu^n\text{-a.e. } x.$$

Using (4.8) in the third term on the left-hand side of (4.6) gives

$$\begin{aligned}
(4.9) \quad & \left\langle \sigma \nabla \log \frac{d\mu^{n+1}}{d\Phi[\mu^n]} \left(T_{\mu^n}^{\mu^{n+1}} \right), T_{\mu^n}^{\mu^{n+1}} - I_d \right\rangle_{L_{\mu^n}^2(\mathbb{R}^d)} = -\tau\sigma^2 \left\| \nabla \log \frac{d\mu^{n+1}}{d\Phi[\mu^n]} \left(T_{\mu^n}^{\mu^{n+1}} \right) \right\|_{L_{\mu^n}^2(\mathbb{R}^d)}^2 \\
& = -\tau\sigma^2 I \left(\mu^{n+1} \middle| \Phi[\mu^n] \right).
\end{aligned}$$

Also, squaring both sides of (4.7) and integrating with respect to μ^{n+1} gives

$$(4.10) \quad I \left(\mu^{n+1} \middle| \Phi[\mu^n] \right) = \frac{1}{\tau^2\sigma^2} \mathcal{W}_2^2(\mu^{n+1}, \mu^n).$$

Using (4.9) and (4.10) in (4.6) gives

$$\begin{aligned}
(4.11) \quad F^\sigma(\mu^{n+1}) &\leq F^\sigma(\mu^n) - \tau\sigma^2 I\left(\mu^{n+1} \middle| \Phi[\mu^n]\right) + 2L'_F \mathcal{W}_2^2(\mu^{n+1}, \mu^n) \\
&= F^\sigma(\mu^n) - \tau\sigma^2 I\left(\mu^{n+1} \middle| \Phi[\mu^n]\right) + 2\tau^2\sigma^2 L'_F I\left(\mu^{n+1} \middle| \Phi[\mu^n]\right) \\
&= F^\sigma(\mu^n) - \tau\sigma^2 (1 - 2\tau L'_F) I\left(\mu^{n+1} \middle| \Phi[\mu^n]\right).
\end{aligned}$$

Now, using again (3.2) with $\mu' = \mu = \mu^{n+1}$, observe that (4.8) can be equivalently written as

$$\begin{aligned}
\frac{1}{\tau} \left(x - T_{\mu^n}^{\mu^{n+1}}(x)\right) &= \nabla_{\mu} F(\mu^n) \left(T_{\mu^n}^{\mu^{n+1}}(x)\right) + \sigma \nabla \log \frac{d\mu^{n+1}}{d\pi} \left(T_{\mu^n}^{\mu^{n+1}}(x)\right) \\
&= \nabla_{\mu} F(\mu^n) \left(T_{\mu^n}^{\mu^{n+1}}(x)\right) - \nabla_{\mu} F(\mu^{n+1}) \left(T_{\mu^n}^{\mu^{n+1}}(x)\right) + \nabla_{\mu} F(\mu^{n+1}) \left(T_{\mu^n}^{\mu^{n+1}}(x)\right) \\
&\quad + \sigma \nabla \log \frac{d\mu^{n+1}}{d\pi} \left(T_{\mu^n}^{\mu^{n+1}}(x)\right) \\
&= \nabla_{\mu} F(\mu^n) \left(T_{\mu^n}^{\mu^{n+1}}(x)\right) - \nabla_{\mu} F(\mu^{n+1}) \left(T_{\mu^n}^{\mu^{n+1}}(x)\right) + \sigma \nabla \log \frac{d\mu^{n+1}}{d\Phi[\mu^{n+1}]} \left(T_{\mu^n}^{\mu^{n+1}}(x)\right).
\end{aligned}$$

By Minkowski's inequality, we obtain

$$\begin{aligned}
\sigma^2 I\left(\mu^{n+1} \middle| \Phi[\mu^{n+1}]\right) &= \sigma^2 \left\| \nabla \log \frac{d\mu^{n+1}}{d\Phi[\mu^{n+1}]} \left(T_{\mu^n}^{\mu^{n+1}}\right) \right\|_{L_{\mu^n}^2(\mathbb{R}^d)}^2 \\
&\leq 2 \left(\left\| \nabla_{\mu} F(\mu^{n+1}) \left(T_{\mu^n}^{\mu^{n+1}}(x)\right) - \nabla_{\mu} F(\mu^n) \left(T_{\mu^n}^{\mu^{n+1}}(x)\right) \right\|_{L_{\mu^n}^2(\mathbb{R}^d)}^2 + \frac{1}{\tau^2} \left\| I_d - T_{\mu^n}^{\mu^{n+1}} \right\|_{L_{\mu^n}^2(\mathbb{R}^d)}^2 \right) \\
&\leq 2(L'_F)^2 \mathcal{W}_2^2(\mu^{n+1}, \mu^n) + \frac{2}{\tau^2} \mathcal{W}_2^2(\mu^{n+1}, \mu^n) \\
&= 2 \left((L'_F)^2 + \frac{1}{\tau^2} \right) \tau^2 \sigma^2 I\left(\mu^{n+1} \middle| \Phi[\mu^n]\right),
\end{aligned}$$

where the second inequality follows from Assumption 2.9 and Corollary A.3, while the last equality follows from (4.10). Hence,

$$I\left(\mu^{n+1} \middle| \Phi[\mu^{n+1}]\right) \leq 2(1 + \tau^2(L'_F)^2) I\left(\mu^{n+1} \middle| \Phi[\mu^n]\right).$$

Since $\tau < \frac{2}{L'_F}$, using this inequality in (4.11) gives

$$\begin{aligned}
F^\sigma(\mu^{n+1}) &\leq F^\sigma(\mu^n) - \frac{\tau\sigma^2(1 - 2\tau L'_F)}{2(1 + \tau^2(L'_F)^2)} I\left(\mu^{n+1} \middle| \Phi[\mu^{n+1}]\right) \\
&\leq F^\sigma(\mu^n) - \frac{\tau\sigma^2(1 - 2\tau L'_F)}{2(1 + \tau^2(L'_F)^2)} \frac{2\alpha_U}{e^{\frac{4C_F}{\sigma}}} \text{KL}\left(\mu^{n+1} \middle| \Phi[\mu^{n+1}]\right) \\
&\leq F^\sigma(\mu^n) - \frac{\tau\sigma(1 - 2\tau L'_F)}{1 + \tau^2(L'_F)^2} \frac{\alpha_U}{e^{\frac{4C_F}{\sigma}}} (F^\sigma(\mu^{n+1}) - F^\sigma(\mu_\sigma^*)),
\end{aligned}$$

where the second and third inequalities follow from (2.7) and Lemma 2.7, respectively.

Let $\kappa := 1 + \frac{\tau\sigma\alpha_U(1-2\tau L'_F)}{(1+\tau^2(L'_F)^2)e^{\frac{4C_F}{\sigma}}} > 1$. Then rearranging the inequality above gives

$$F^\sigma(\mu^{n+1}) - F^\sigma(\mu_\sigma^*) \leq \kappa^{-1} (F^\sigma(\mu^n) - F^\sigma(\mu_\sigma^*)).$$

The convergence estimate follows by iterating over $n \in \mathbb{N}$.

4.3. Proof of (iii). First, note that by Proposition 7.4, the iterates $(\mu^n)_{n \in \mathbb{N}} \subset \mathfrak{C}$, and by Lemma 7.1 that $(\nu^n)_{n \in \mathbb{N}} \subset \mathcal{P}_2^\lambda(\mathbb{R}^d)$. By Proposition 7.5, we have

$$\sigma \nabla \log \frac{d\mu^{n+1}}{d\pi}(x) = \frac{1}{\tau} \left(T_{\nu^{n+1}}^{\mu^{n+1}}(x) - x \right), \quad \text{for } \mu^{n+1}\text{-a.e. } x.$$

By Corollary 7.3, exists a unique ν^{n+1} -a.e. optimal transport map $T_{\nu^{n+1}}^{\mu^{n+1}} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that $T_{\mu^{n+1}}^{\nu^{n+1}} \circ T_{\nu^{n+1}}^{\mu^{n+1}} = I_d$, ν^{n+1} -a.e.. We thus have

$$(4.12) \quad T_{\nu^{n+1}}^{\mu^{n+1}} = I_d - \tau \sigma \nabla \log \frac{d\mu^{n+1}}{d\pi} \left(T_{\nu^{n+1}}^{\mu^{n+1}} \right), \quad \nu^{n+1}\text{-a.e..}$$

By convexity of $\text{KL}(\cdot|\pi)$ along generalized geodesics in the Wasserstein space [1, Theorem 9.4.10], taking $\mu = \mu^{n+1}$, $\pi = \mu^n$ and $\nu = \nu^{n+1}$ in [31, Lemma 4] gives

$$\text{KL}(\mu^{n+1}|\pi) \leq \text{KL}(\mu^n|\pi) - \left\langle \nabla \log \frac{d\mu^{n+1}}{d\pi} \left(T_{\nu^{n+1}}^{\mu^{n+1}} \right), T_{\nu^{n+1}}^{\mu^n} - T_{\nu^{n+1}}^{\mu^{n+1}} \right\rangle_{L^2_{\nu^{n+1}}(\mathbb{R}^d)}.$$

By Corollary 7.3, we have $T_{\nu^{n+1}}^{\mu^n} = (I_d - \tau \nabla_{\mu} F(\mu^n)(\cdot))^{-1}$, ν^{n+1} -a.e.. Therefore,

$$\text{KL}(\mu^{n+1}|\pi) \leq \text{KL}(\mu^n|\pi) - \left\langle \nabla \log \frac{d\mu^{n+1}}{d\pi} \left(T_{\nu^{n+1}}^{\mu^{n+1}} \right), (I_d - \tau \nabla_{\mu} F(\mu^n))^{-1} - T_{\nu^{n+1}}^{\mu^{n+1}} \right\rangle_{L^2_{\nu^{n+1}}(\mathbb{R}^d)}.$$

Let us denote the pushforward from μ^n to μ^{n+1} by $P_{\mu^n}^{\mu^{n+1}} := T_{\nu^{n+1}}^{\mu^{n+1}} \circ T_{\mu^n}^{\nu^{n+1}}$. Then $T_{\nu^{n+1}}^{\mu^{n+1}} = P_{\mu^n}^{\mu^{n+1}} \circ (I_d - \tau \nabla_{\mu} F(\mu^n))^{-1}$, and hence by A.1, the last inequality is equivalent to

$$(4.13) \quad \text{KL}(\mu^{n+1}|\pi) \leq \text{KL}(\mu^n|\pi) - \left\langle \nabla \log \frac{d\mu^{n+1}}{d\pi} \left(P_{\mu^n}^{\mu^{n+1}} \right), I_d - P_{\mu^n}^{\mu^{n+1}} \right\rangle_{L^2_{\mu^n}(\mathbb{R}^d)}.$$

Using (4.12) and $P_{\mu^n}^{\mu^{n+1}} = T_{\nu^{n+1}}^{\mu^{n+1}} \circ (I_d - \tau \nabla_{\mu} F(\mu^n))$, we have

$$(4.14) \quad P_{\mu^n}^{\mu^{n+1}} = I_d - \tau \nabla_{\mu} F(\mu^n) - \tau \sigma \nabla \log \frac{d\mu^{n+1}}{d\pi} \left(P_{\mu^n}^{\mu^{n+1}} \right), \quad \mu^n\text{-a.e.,}$$

Note that, for each $n \in \mathbb{N}$, $\gamma^n := \left(I_d, P_{\mu^n}^{\mu^{n+1}} \right)_{\#} \mu^n$ is a coupling between μ^n and μ^{n+1} . Then, by Minkowski's inequality, we obtain

$$(4.15) \quad \begin{aligned} \sigma^2 I \left(\mu^{n+1} \middle| \Phi[\mu^{n+1}] \right) &= \sigma^2 \left\| \nabla \log \frac{d\mu^{n+1}}{d\Phi[\mu^{n+1}]} \left(P_{\mu^n}^{\mu^{n+1}} \right) \right\|_{L^2_{\mu^n}(\mathbb{R}^d)}^2 \\ &= \left\| \nabla_{\mu} F(\mu^{n+1}) \left(P_{\mu^n}^{\mu^{n+1}} \right) + \sigma \log \frac{d\mu^{n+1}}{d\pi} \left(P_{\mu^n}^{\mu^{n+1}} \right) \right\|_{L^2_{\mu^n}(\mathbb{R}^d)}^2 \\ &\leq 2 \left(\left\| \nabla_{\mu} F(\mu^{n+1}) \left(P_{\mu^n}^{\mu^{n+1}} \right) - \nabla_{\mu} F(\mu^n) \right\|_{L^2_{\mu^n}(\mathbb{R}^d)}^2 + \frac{1}{\tau^2} \left\| I_d - P_{\mu^n}^{\mu^{n+1}} \right\|_{L^2_{\mu^n}(\mathbb{R}^d)}^2 \right) \\ &\leq 2 \left(2(L'_F)^2 \mathcal{W}_2^2(\mu^{n+1}, \mu^n) + 2(L'_F)^2 \left\| I_d - P_{\mu^n}^{\mu^{n+1}} \right\|_{L^2_{\mu^n}(\mathbb{R}^d)}^2 + \frac{1}{\tau^2} \left\| I_d - P_{\mu^n}^{\mu^{n+1}} \right\|_{L^2_{\mu^n}(\mathbb{R}^d)}^2 \right) \\ &\leq 2 \left(\frac{1}{\tau^2} + 4(L'_F)^2 \right) \left\| I_d - P_{\mu^n}^{\mu^{n+1}} \right\|_{L^2_{\mu^n}(\mathbb{R}^d)}^2, \end{aligned}$$

where the first equality follows from (3.2) with $\mu' = \mu = \mu^{n+1}$, first inequality follows from (4.14) and last two inequalities follow from Assumption 2.9 and (A.2), respectively.

Multiplying (4.13) by σ and using Lemma 6.5 gives

$$\begin{aligned} F^\sigma(\mu^{n+1}) &\leq F^\sigma(\mu^n) + \left\langle \nabla_\mu F(\mu^n)(\cdot) + \sigma \nabla \log \frac{d\mu^{n+1}}{d\pi} \left(P_{\mu^n}^{\mu^{n+1}} \right), P_{\mu^n}^{\mu^{n+1}} - I_d \right\rangle_{L_{\mu^n}^2(\mathbb{R}^d)} \\ &\quad + L'_F \left\| I_d - P_{\mu^n}^{\mu^{n+1}} \right\|_{L_{\mu^n}^2(\mathbb{R}^d)}^2 \\ &= F^\sigma(\mu^n) - \frac{1}{\tau} \left\| I_d - P_{\mu^n}^{\mu^{n+1}} \right\|_{L_{\mu^n}^2(\mathbb{R}^d)}^2 + L'_F \left\| I_d - P_{\mu^n}^{\mu^{n+1}} \right\|_{L_{\mu^n}^2(\mathbb{R}^d)}^2, \end{aligned}$$

where the equality follows from (4.14). Hence

$$\begin{aligned} F^\sigma(\mu^{n+1}) &\leq F^\sigma(\mu^n) - \left(\frac{1}{\tau} - L'_F \right) \left\| I_d - P_{\mu^n}^{\mu^{n+1}} \right\|_{L_{\mu^n}^2(\mathbb{R}^d)}^2 \\ &\leq F^\sigma(\mu^n) - \left(\frac{1}{\tau} - L'_F \right) \frac{\sigma^2}{2 \left(\frac{1}{\tau^2} + 4(L'_F)^2 \right)} I \left(\mu^{n+1} \middle| \Phi[\mu^{n+1}] \right) \\ &= F^\sigma(\mu^n) - (1 - \tau L'_F) \frac{\tau \sigma^2}{2(1 + 4(\tau L'_F)^2)} I \left(\mu^{n+1} \middle| \Phi[\mu^{n+1}] \right) \\ &\leq F^\sigma(\mu^n) - (1 - \tau L'_F) \frac{\tau \alpha_U \sigma^2}{e^{\frac{4C_F}{\sigma}} (1 + 4(\tau L'_F)^2)} \text{KL} \left(\mu^{n+1} \middle| \Phi[\mu^{n+1}] \right), \end{aligned}$$

the first inequality follows from (4.15) and the fact that $\tau < \frac{1}{L'_F}$, whereas the last inequality follows from (2.7). Let

$$\kappa := 1 + (1 - \tau L'_F) \frac{\tau \alpha_U \sigma}{e^{\frac{4C_F}{\sigma}} (1 + 4(\tau L'_F)^2)} > 1.$$

Then using Lemma 2.7 in the previous inequality gives

$$F^\sigma(\mu^{n+1}) - F^\sigma(\mu_\sigma^*) \leq \kappa^{-1} (F^\sigma(\mu^n) - F^\sigma(\mu_\sigma^*)).$$

The convergence estimate follows by iterating over $n \in \mathbb{N}$.

5. PROXIMAL POINT SCHEME

In this section, we present the auxiliary results needed for the proof of (i) in Theorem 3.1. We start by proving that (1.3) admits a unique minimizer.

Theorem 5.1 (Existence and uniqueness of minimizer for (1.3)). *Let $F \in \mathcal{C}^1$ and Assumption 2.1, 2.2 hold. Given $\mu^0 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, there exists a unique minimizer $\mu^1 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$ of*

$$(5.1) \quad \mathcal{P}_2^\pi(\mathbb{R}^d) \ni \mu \mapsto \mathcal{F}(\mu) := F^\sigma(\mu) + \frac{1}{2\tau} \mathcal{W}_2^2(\mu, \mu^0).$$

The proof is a modification of the argument in the proof of [20, Proposition 4.1], and before we present it we give an outline of the main steps:

Step 1. Firstly, we show that \mathcal{F} is bounded below on $\mathcal{P}_2^\pi(\mathbb{R}^d)$.

Step 2. Secondly, we show that any minimizing sequence $(\mu_k)_{k \in \mathbb{N}} \subset \mathcal{P}_2^\pi(\mathbb{R}^d)$ of \mathcal{F} contains a weakly convergent subsequence in $L_\pi^1(\mathbb{R}^d)$.

Step 3. Thirdly, we show that the weak limit of the converging subsequence is indeed a minimizer of \mathcal{F} .

Step 4. Finally, we deduce the uniqueness of the minimizer from the strict convexity of $\text{KL}(\cdot | \pi)$.

Proof. Step 1. By Jensen's inequality since the map $z \mapsto z \log z$ is convex on $[0, \infty)$, it follows that

$$(5.2) \quad \text{KL}(\mu|\pi) \geq 0, \quad \text{for all } \mu \in \mathcal{P}_2^\pi(\mathbb{R}^d).$$

Since $\frac{1}{2\tau}\mathcal{W}_2^2(\cdot, \mu^0) \geq 0$ and the fact that F is bounded below on $\mathcal{P}_2^\pi(\mathbb{R}^d)$, by Assumption 2.2, it follows that \mathcal{F} is bounded below on $\mathcal{P}_2^\pi(\mathbb{R}^d)$, and thus $\inf_{\mu \in \mathcal{P}_2^\pi(\mathbb{R}^d)} \mathcal{F}(\mu) > -\infty$.

Step 2. Let $(\mu_k)_{k \in \mathbb{N}} \subset \mathcal{P}_2^\pi(\mathbb{R}^d)$ be a minimizing sequence for \mathcal{F} , i.e., $\lim_{k \rightarrow \infty} \mathcal{F}(\mu_k) = \inf_{\mu \in \mathcal{P}_2^\pi(\mathbb{R}^d)} \mathcal{F}(\mu)$. Then the sequence $(\mathcal{F}(\mu_k))_k$ is bounded on $\mathcal{P}_2^\pi(\mathbb{R}^d)$, i.e., there exists $M_{\mathcal{F}} > 0$ such that $|\mathcal{F}(\mu_k)| \leq M_{\mathcal{F}}$, for all $k \in \mathbb{N}$.

Thus, since $\frac{1}{2\tau}\mathcal{W}_2^2(\cdot, \mu^0) \geq 0$, it follows that

$$\text{KL}(\mu_k|\pi) \leq \frac{1}{\sigma} (M_{\mathcal{F}} - F(\mu_k)) < \frac{1}{\sigma} \left(M_{\mathcal{F}} - \inf_{\mu \in \mathcal{P}_2^\pi(\mathbb{R}^d)} F(\mu) \right) < \infty,$$

and together with (5.2),

$$(5.3) \quad (\text{KL}(\mu_k|\pi))_k \text{ is bounded.}$$

From the inequality $|y|^2 \leq 2|x|^2 + 2|x - y|^2$, which holds for all $x, y \in \mathbb{R}^d$, and (A.2), it follows that

$$(5.4) \quad \int_{\mathbb{R}^d} |y|^2 \mu''(\text{d}y) \leq \int_{\mathbb{R}^d} |x|^2 \mu'(\text{d}x) + 2\mathcal{W}_2^2(\mu', \mu''), \quad \text{for all } \mu', \mu'' \in \mathcal{P}_2^\pi(\mathbb{R}^d),$$

Again using (5.2) we obtain

$$M_{\mathcal{F}} \geq \mathcal{F}(\mu_k) \geq F(\mu_k) + \frac{1}{2\tau}\mathcal{W}_2^2(\mu_k, \mu^0) \geq F(\mu_k) + \frac{1}{4\tau} \int_{\mathbb{R}^d} |x|^2 \mu_k(\text{d}x) - \frac{1}{2\tau} \int_{\mathbb{R}^d} |x|^2 \mu^0(\text{d}x).$$

Hence, by Assumption 2.2,

$$\begin{aligned} \int_{\mathbb{R}^d} |x|^2 \mu_k(\text{d}x) &\leq 4\tau (M_{\mathcal{F}} - F(\mu_k)) + 2 \int_{\mathbb{R}^d} |x|^2 \mu^0(\text{d}x) \\ &\leq 4\tau \left(M_{\mathcal{F}} - \inf_{\mu \in \mathcal{P}_2^\pi(\mathbb{R}^d)} F(\mu) \right) + 2 \int_{\mathbb{R}^d} |x|^2 \mu^0(\text{d}x) < \infty. \end{aligned}$$

and hence

$$(5.5) \quad \left(\int_{\mathbb{R}^d} |x|^2 \mu_k(\text{d}x) \right)_k \text{ is bounded.}$$

Note that there exists $C > 0$ such that

$$|\min\{z \log z, 0\}| \leq C, \quad \text{for all } z \geq 0.$$

Hence, we obtain

$$(5.6) \quad \int_{\mathbb{R}^d} \left| \min \left\{ \frac{\text{d}\mu}{\text{d}\pi}(x) \log \frac{\text{d}\mu}{\text{d}\pi}(x), 0 \right\} \right| \pi(\text{d}x) \leq C \int_{\mathbb{R}^d} \pi(\text{d}x) = C.$$

Furthermore, from (5.6), we obtain

$$(5.7) \quad \left(\int_{\mathbb{R}^d} \left| \min \left\{ \frac{\text{d}\mu_k}{\text{d}\pi}(x) \log \frac{\text{d}\mu_k}{\text{d}\pi}(x), 0 \right\} \right| \pi(\text{d}x) \right)_k \text{ is bounded.}$$

Since $\max\{z \log z, 0\} = z \log z + |\min\{z \log z, 0\}|$, for all $z \geq 0$, it follows from (5.3) and (5.7) that

$$\left(\int_{\mathbb{R}^d} \max \left\{ \frac{\text{d}\mu_k}{\text{d}\pi}(x) \log \frac{\text{d}\mu_k}{\text{d}\pi}(x), 0 \right\} \pi(\text{d}x) \right)_k \text{ is bounded.}$$

Since $\|\frac{d\mu_k}{d\pi}\|_{L^1_\pi(\mathbb{R}^d)} = 1$, for all $k \in \mathbb{N}$, we obtain that $(\frac{d\mu_k}{d\pi})_k$ is uniformly bounded in $L^1_\pi(\mathbb{R}^d)$. As $[0, \infty) \ni z \mapsto \max\{z \log z, 0\}$ is non-negative, increasing and has superlinear growth together with (5.3) implies via [5, Theorem 4.5.9] that $(\frac{d\mu_k}{d\pi})_k$ is uniformly integrable. Consequently, according to the Dunford-Pettis theorem (see [5, Corollary 4.7.19]), there exists $\mu^* \in \mathcal{P}_2^\pi(\mathbb{R}^d)$ such that (at least for a subsequence)

$$(5.8) \quad \frac{d\mu_k}{d\pi} \rightarrow \frac{d\mu^*}{d\pi} \text{ weakly in } L^1_\pi(\mathbb{R}^d) \text{ as } k \rightarrow \infty,$$

i.e.,

$$\int_{\mathbb{R}^d} h(x) \frac{d\mu_k}{d\pi}(x) \pi(dx) \rightarrow \int_{\mathbb{R}^d} h(x) \frac{d\mu^*}{d\pi}(x) \pi(dx) \text{ as } k \rightarrow \infty,$$

for all $h \in L^\infty(\mathbb{R}^d)$.

Step 3. Observe that any continuous bounded function $g : \mathbb{R}^d \rightarrow \mathbb{R}$ is in $L^\infty_\pi(\mathbb{R}^d)$, and hence, by (5.8), we obtain as $k \rightarrow \infty$,

$$\int_{\mathbb{R}^d} g(x) \mu_k(dx) = \int_{\mathbb{R}^d} g(x) \frac{d\mu_k}{d\pi}(x) \pi(dx) \rightarrow \int_{\mathbb{R}^d} g(x) \frac{d\mu^*}{d\pi}(x) \pi(dx) = \int_{\mathbb{R}^d} g(x) \mu^*(dx),$$

i.e., $\mu_k \rightarrow \mu^*$ weakly (with respect to the topology of probability measures convergence) as $k \rightarrow \infty$.

Since $\text{KL}(\cdot|\pi)$ is lower semi-continuous with respect to weak convergence of probability measures, it follows that

$$\text{KL}(\mu^*|\pi) \leq \liminf_{k \rightarrow \infty} \text{KL}(\mu_k|\pi).$$

By (5.5) and continuity of \mathcal{W}_2 (see [36, Corollary 6.11]), we have

$$\mathcal{W}_2^2(\mu^0, \mu^*) = \lim_{k \rightarrow \infty} \mathcal{W}_2^2(\mu^0, \mu_k).$$

Now, we show that $\lim_{k \rightarrow \infty} F(\mu_k) = F(\mu^*)$. Indeed, by Definition B.1, we have

$$|F(\mu_k) - F(\mu^*)| \leq \int_0^1 \left| \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu_{\lambda,k}, x) (\mu_k - \mu)(dx) \right| d\eta,$$

where $\mu_{\eta,k} := (1 - \eta)\mu_k + \eta\mu^*$. For every $\eta \in [0, 1]$, we have

$$\begin{aligned} & \left| \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu_{\eta,k}, x) (\mu_k - \mu^*)(dx) \right| \\ & \leq \left| \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu^*, x) (\mu_k - \mu^*)(dx) \right| + \int_{\mathbb{R}^d} \left| \frac{\delta F}{\delta \mu}(\mu_{\eta,k}, x) - \frac{\delta F}{\delta \mu}(\mu^*, x) \right| (\mu_k + \mu^*)(dx). \end{aligned}$$

Since $\frac{\delta F}{\delta \mu}(\mu^*, \cdot)$ is a bounded continuous function (cf. Definition B.1 and Assumption 2.2), the weak convergence $\mu_k \rightarrow \mu^*$ (with respect to convergence of probability measures) implies

$$\lim_{k \rightarrow \infty} \left| \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu^*, x) (\mu_k - \mu^*)(dx) \right| = 0.$$

For the second term, by Assumption 2.2, we have

$$\limsup_{k \rightarrow \infty} \int_{\mathbb{R}^d} \left| \frac{\delta F}{\delta \mu}(\mu_{\eta,k}, x) - \frac{\delta F}{\delta \mu}(\mu^*, x) \right| (\mu_k + \mu^*)(dx) \leq 2L_F \limsup_{k \rightarrow \infty} \mathcal{W}_2(\mu_{\eta,k}, \mu^*) = 0,$$

by (5.5) and continuity of \mathcal{W}_2 (see [36, Corollary 6.11]). Finally, using Assumption 2.2, we have, for any $\eta \in [0, 1]$,

$$\left| \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu_{\eta,k}, x) (\mu_k - \mu^*)(dx) \right| \leq 2C_F,$$

thus we can apply the dominated convergence theorem, and obtain that as $k \rightarrow \infty$,

$$\int_0^1 \left| \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu_{\eta,k}, x)(\mu_k - \mu^*)(dx) \right| d\eta \rightarrow 0.$$

Putting everything together, we obtain

$$\mathcal{F}(\mu^*) \leq \liminf_{k \rightarrow \infty} \mathcal{F}(\mu_k) = \inf_{\mu \in \mathcal{P}_2^\pi(\mathbb{R}^d)} \mathcal{F}(\mu).$$

On the other hand, from the definition of infimum, we have

$$\mathcal{F}(\mu^*) \geq \inf_{\mu \in \mathcal{P}_2^\pi(\mathbb{R}^d)} \mathcal{F}(\mu).$$

Hence, $\mathcal{F}(\mu^*) = \inf_{\mu \in \mathcal{P}_2^\pi(\mathbb{R}^d)} \mathcal{F}(\mu)$, and therefore $\mu^* \in \mathcal{P}_2^\pi(\mathbb{R}^d)$ is a minimizer of \mathcal{F} , which we shall denote by μ^1 .

Step 4. The uniqueness of the minimizer of \mathcal{F} follows from Assumption 2.1, convexity of $\mathcal{P}_2^\pi(\mathbb{R}^d) \ni \mu \mapsto \mathcal{W}_2^2(\mu, \mu^0)$, and strict convexity of $\mathcal{P}_2^\pi(\mathbb{R}^d) \ni \mu \mapsto \text{KL}(\mu|\pi)$. \square

From Theorem 5.1 it follows inductively that, for each $n \in \mathbb{N}$, given $\mu^n \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, the scheme (1.3) admits a unique minimizer $\mu^{n+1} \in \mathcal{P}_2^\pi(\mathbb{R}^d)$. Hence, if $\mu^0 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, then $(\mu^n)_{n \in \mathbb{N}} \subset \mathcal{P}_2^\pi(\mathbb{R}^d)$ along the scheme (1.3). Therefore, via Theorem A.2, we obtain

Corollary 5.2 (Existence of optimal transport maps along (1.3)). *Let $\nu \in \mathcal{P}_2(\mathbb{R}^d)$. Let $F \in \mathcal{C}^1$ and Assumption 2.1, 2.2 hold. Given $\mu^0 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, there exists a unique μ^n -a.e. optimal transport map $T_{\mu^n}^\nu : \mathbb{R}^d \rightarrow \mathbb{R}^d$ from μ^n to ν . In particular, if $\nu = \mu^{n+1}$, there also exists a unique μ^{n+1} -a.e. optimal transport map $T_{\mu^{n+1}}^{\mu^n} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that $T_{\mu^{n+1}}^{\mu^n} \circ T_{\mu^n}^{\mu^{n+1}} = I_d$, μ^n -a.e and $T_{\mu^n}^{\mu^{n+1}} \circ T_{\mu^{n+1}}^{\mu^n} = I_d$, μ^{n+1} -a.e..*

The following proposition together with an induction argument guarantee that $\text{KL}(\cdot|\pi)$ admits a unique Wasserstein sub-differential given by $\nabla \log \frac{d\mu^n}{d\pi}$, where μ^n is the iterate generated by (1.3) at each step $n \geq 1$.

Proposition 5.3 (Wasserstein sub-differentiability class for $\text{KL}(\cdot|\pi)$ along (1.3)). *Let (2.4) in Assumption 2.3 hold. Let $F \in \mathcal{C}^1$ and Assumption 2.1, 2.2 hold. Given $\mu^0 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, the unique minimizer $\mu^1 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$ of (5.1) belongs to \mathfrak{C} .*

Proof. Since $\mu^0 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, $F \in \mathcal{C}^1$, and Assumption 2.1, 2.2 hold, Theorem 5.1 guarantees the existence and uniqueness of a minimizer $\mu^1 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$ for (5.1). We will now prove that $\mu^1 \in \mathfrak{C}$. By [1, Definition 1.2.4], the metric slope of the relative entropy $\text{KL}(\cdot|\pi)$ at $\mu \in \mathcal{P}_2^\pi(\mathbb{R}^d)$ is defined by

$$|\mathfrak{D} \text{KL}(\cdot|\pi)|(\mu) := \limsup_{\nu \rightarrow \mu} \frac{(\text{KL}(\mu|\pi) - \text{KL}(\nu|\pi))_+}{\mathcal{W}_2(\nu, \mu)}.$$

Since $\mu^1 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$ is a minimizer for (5.1), it follows that for any $\mu \in \mathcal{P}_2^\pi(\mathbb{R}^d)$,

$$\begin{aligned} \text{KL}(\mu^1|\pi) - \text{KL}(\mu|\pi) &\leq \frac{1}{\sigma} (F(\mu) - F(\mu^1)) + \frac{1}{2\tau\sigma} (\mathcal{W}_2^2(\mu, \mu^0) - \mathcal{W}_2^2(\mu^1, \mu^0)) \\ (5.9) \quad &= \frac{1}{\sigma} \int_0^1 \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu_\eta, x)(\mu - \mu^1)(dx) d\eta \\ &+ \frac{1}{2\tau\sigma} (\mathcal{W}_2(\mu, \mu^0) - \mathcal{W}_2(\mu^1, \mu^0)) (\mathcal{W}_2(\mu, \mu^0) + \mathcal{W}_2(\mu^1, \mu^0)) \\ &\leq \frac{1}{\sigma} \int_0^1 \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu_\eta, x)(\mu - \mu^1)(dx) d\eta + \frac{1}{2\tau\sigma} \mathcal{W}_2(\mu, \mu^1) (\mathcal{W}_2(\mu, \mu^0) + \mathcal{W}_2(\mu^1, \mu^0)), \end{aligned}$$

where the first equality follows from Definition B.1 with $\mu_\eta := \mu + \eta(\mu^1 - \mu)$, and the last inequality follows from the triangle inequality applied to \mathcal{W}_2 .

By Theorem A.2, there exists a unique optimal coupling $\gamma^* \in \Gamma_o(\mu, \mu^1)$, and hence

$$\begin{aligned} \left| \int_0^1 \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu_\eta, x)(\mu - \mu^1)(dx) d\eta \right| &= \left| \int_0^1 \int_{\mathbb{R}^d \times \mathbb{R}^d} \left(\frac{\delta F}{\delta \mu}(\mu_\eta, x) - \frac{\delta F}{\delta \mu}(\mu_\eta, y) \right) \gamma^*(dx, dy) d\eta \right| \\ &\leq \int_0^1 \int_{\mathbb{R}^d \times \mathbb{R}^d} \left| \frac{\delta F}{\delta \mu}(\mu_\eta, x) - \frac{\delta F}{\delta \mu}(\mu_\eta, y) \right| \gamma^*(dx, dy) d\eta \\ &\leq \left(\int_0^1 \int_{\mathbb{R}^d \times \mathbb{R}^d} \left| \frac{\delta F}{\delta \mu}(\mu_\eta, x) - \frac{\delta F}{\delta \mu}(\mu_\eta, y) \right|^2 \gamma^*(dx, dy) d\eta \right)^{1/2} \\ &\leq L_F \left(\int_{\mathbb{R}^d \times \mathbb{R}^d} |x - y|^2 \gamma^*(dx, dy) \right)^{1/2} = L_F \mathcal{W}_2(\mu, \mu^1), \end{aligned}$$

where the penultimate inequality follows from the Cauchy-Schwarz inequality and the fact that $\gamma^* \in \mathcal{P}_2(\mathbb{R}^d \times \mathbb{R}^d)$, the last inequality follows from Assumption 2.2 and the last equality follows from optimality of γ^* .

Therefore, for any $\mu \neq \mu^1$, dividing (5.9) by $\mathcal{W}_2(\mu, \mu^1)$ gives

$$\frac{(\text{KL}(\mu^1|\pi) - \text{KL}(\mu|\pi))_+}{\mathcal{W}_2(\mu, \mu^1)} \leq \frac{L_F}{\sigma} + \frac{1}{2\tau\sigma} (\mathcal{W}_2(\mu, \mu^0) + \mathcal{W}_2(\mu^1, \mu^0)).$$

Taking limsup as $\mu \rightarrow \mu^1$ and using continuity of \mathcal{W}_2 via [36, Corollary 6.11] yields

$$|\mathfrak{D} \text{KL}(\cdot|\pi)|(\mu^1) \leq \frac{L_F}{\sigma} + \frac{1}{\tau\sigma} \mathcal{W}_2(\mu^1, \mu^0) < \infty.$$

Hence, using Theorem A.5, it follows that $\frac{d\mu^1}{d\pi} \in W_{\lambda, \text{loc}}^{1,1}(\mathbb{R}^d)$, $\frac{|\nabla \frac{d\mu^1}{d\pi}|^2}{\frac{d\mu^1}{d\pi}} \in L^1_\pi(\mathbb{R}^d)$, and

$|\mathfrak{D} \text{KL}(\cdot|\pi)|^2(\mu^1) = I(\mu^1|\pi)$. Since $\frac{|\nabla \frac{d\mu^1}{d\pi}|^2}{\frac{d\mu^1}{d\pi}} \in L^1_\pi(\mathbb{R}^d)$ we have $\nabla \sqrt{\frac{d\mu^1}{d\pi}} \in L^2_\pi(\mathbb{R}^d)$. In addition, $\mu^1 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$ implies $\sqrt{\frac{d\mu^1}{d\pi}} \in L^2_\pi(\mathbb{R}^d)$, and therefore we obtain that $\mu^1 \in \mathfrak{C}$. \square

From Proposition 5.3 it follows inductively that, for each $n \in \mathbb{N}$, given $\mu^n \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, the unique minimizer μ^{n+1} of (1.3) belongs to \mathfrak{C} . Hence, if $\mu^0 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, then $(\mu^n)_{n \in \mathbb{N}} \subset \mathfrak{C}$ along the scheme (1.3).

Since $\text{KL}(\cdot|\pi)$ admits a unique Wasserstein sub-differential and F is Wasserstein differentiable by Assumption 2.8, [1, Lemma 10.1.2] allows us to write the following first-order optimality condition for (1.3).

Lemma 5.4 (Optimality condition for (1.3)). *Let $F \in \mathcal{C}^1$, Assumption 2.1, 2.2, 2.8 and (2.4) in Assumption 2.3 hold. Given $\mu^0 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, then, for each $n \in \mathbb{N}$, the unique minimizer $\mu^{n+1} \in \mathfrak{C}$ of (1.3) satisfies*

$$\nabla_\mu F(\mu^{n+1})(x) + \sigma \nabla \log \frac{d\mu^{n+1}}{d\pi}(x) = \frac{1}{\tau} \left(T_{\mu^{n+1}}^{\mu^n}(x) - x \right), \quad \text{for } \mu^{n+1}\text{-a.e. } x.$$

Proof. For each $n \in \mathbb{N}$, let us denote

$$J_n(\mu) := F(\mu) + \sigma \text{KL}(\mu|\pi),$$

for all $\mu \in \mathcal{P}_2^\pi(\mathbb{R}^d)$. Note that J_n is lower semi-continuous and $J_n < +\infty$. From Proposition 5.3, we have that (1.3) admits a unique minimizer $\mu^{n+1} \in \mathfrak{C}$. By Corollary 5.2, there exists

a unique μ^{n+1} -a.e. optimal transport map $T_{\mu^{n+1}}^{\mu^n} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ from μ^{n+1} to μ^n . Hence, by [1, Lemma 10.1.2],

$$\frac{1}{\tau} \left(T_{\mu^{n+1}}^{\mu^n} - I \right) \in \partial^- J_n(\mu^{n+1}),$$

where ∂^- denotes the Wasserstein sub-differential (cf. Definition B.3). By Assumption 2.8 and using the fact that $\mu^{n+1} \in \mathfrak{C}$ together with Theorem A.5, we obtain

$$\partial^- J_n(\mu^{n+1}) = \left\{ \nabla_{\mu} F(\mu^{n+1}) + \sigma \nabla \log \frac{d\mu^{n+1}}{d\pi} \right\},$$

and hence the conclusion follows. \square

6. PROX-LINEAR SCHEME

In this section, we present the auxiliary results needed for the proof of (ii) in Theorem 3.1. We start by proving that (1.4) admits a unique minimizer.

Theorem 6.1 (Existence and uniqueness of minimizer for (1.4)). *Let $F \in \mathcal{C}^1$ and (2.3) in Assumption 2.2 hold. Given $\mu^0 \in \mathcal{P}_2^{\pi}(\mathbb{R}^d)$, there exists a unique minimizer $\mu^1 \in \mathcal{P}_2^{\pi}(\mathbb{R}^d)$ of*

$$(6.1) \quad \mathcal{P}_2^{\pi}(\mathbb{R}^d) \ni \mu \mapsto \mathcal{G}(\mu) := \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu^0, x)(\mu - \mu^0)(dx) + \sigma \text{KL}(\mu|\pi) + \frac{1}{2\tau} \mathcal{W}_2^2(\mu, \mu^0).$$

Proof. The proof follows the same steps as the proof of Theorem 5.1. Observe that

$$\operatorname{argmin}_{\mu \in \mathcal{P}_2^{\pi}(\mathbb{R}^d)} \mathcal{G}(\mu) = \operatorname{argmin}_{\mu \in \mathcal{P}_2^{\pi}(\mathbb{R}^d)} \left\{ \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu^0, x)\mu(dx) + \sigma \text{KL}(\mu|\pi) + \frac{1}{2\tau} \mathcal{W}_2^2(\mu, \mu^0) \right\},$$

and therefore it suffices to show that $\mu^1 \in \mathcal{P}_2^{\pi}(\mathbb{R}^d)$ is the unique minimizer of the function on the right-hand side, which we denote \mathcal{G} . Note that, for any $\mu \in \mathcal{P}_2^{\pi}(\mathbb{R}^d)$,

$$\begin{aligned} \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu^0, x)\mu(dx) + \sigma \text{KL}(\mu|\pi) &= -\sigma \int_{\mathbb{R}^d} \log e^{-\frac{1}{\sigma} \frac{\delta F}{\delta \mu}(\mu^0, x)} \mu(dx) + \sigma \int_{\mathbb{R}^d} \log \frac{d\mu}{d\pi}(x)\mu(dx) \\ &= \sigma \text{KL}(\mu|\Phi[\mu^0]), \end{aligned}$$

where the last equality follows from (2.6). Hence, using (2.3) in Assumption 2.2, the result follows from [20, Proposition 4.1]. \square

From Theorem 6.1 it follows inductively that, for each $n \in \mathbb{N}$, given $\mu^n \in \mathcal{P}_2^{\pi}(\mathbb{R}^d)$, the scheme (1.4) admits a unique minimizer $\mu^{n+1} \in \mathcal{P}_2^{\pi}(\mathbb{R}^d)$. Hence, if $\mu^0 \in \mathcal{P}_2^{\pi}(\mathbb{R}^d)$, then $(\mu^n)_{n \in \mathbb{N}} \subset \mathcal{P}_2^{\pi}(\mathbb{R}^d)$ along the scheme (1.4). Therefore, via Theorem A.2, we obtain

Corollary 6.2 (Existence of optimal transport maps along (1.4)). *Let $\nu \in \mathcal{P}_2(\mathbb{R}^d)$. Let $F \in \mathcal{C}^1$ and (2.3) in Assumption 2.2 hold. Given $\mu^0 \in \mathcal{P}_2^{\pi}(\mathbb{R}^d)$, there exists a unique μ^n -a.e. optimal transport map $T_{\mu^n}^{\nu} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ from μ^n to ν . In particular, if $\nu = \mu^{n+1}$, there also exists a unique μ^{n+1} -a.e. optimal transport map $T_{\mu^{n+1}}^{\mu^n} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that $T_{\mu^{n+1}}^{\mu^n} \circ T_{\mu^n}^{\mu^{n+1}} = I_d$, μ^n -a.e and $T_{\mu^n}^{\mu^{n+1}} \circ T_{\mu^{n+1}}^{\mu^n} = I_d$, μ^{n+1} -a.e..*

The following proposition together with an induction argument guarantee that $\text{KL}(\cdot|\pi)$ admits a unique Wasserstein sub-differential given by $\nabla \log \frac{d\mu^n}{d\pi}$, where μ^n is the iterate generated by (1.4) at each step $n \geq 1$.

Proposition 6.3 (Wasserstein sub-differentiability class for $\text{KL}(\cdot|\pi)$ along (1.4)). *Let (2.4) in Assumption 2.3 hold. Let $F \in \mathcal{C}^1$ and Assumption 2.2 hold. Given $\mu^0 \in \mathcal{P}_2^{\pi}(\mathbb{R}^d)$, the unique minimizer $\mu^1 \in \mathcal{P}_2^{\pi}(\mathbb{R}^d)$ of (6.1) belongs to \mathfrak{C} .*

Proof. The proof follows the same steps as the proof of Proposition 5.3. Since $\mu^0 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, $F \in \mathcal{C}^1$ and (2.3) in Assumption 2.2 holds, Theorem 6.1 guarantees the existence and uniqueness of a minimizer $\mu^1 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$ for (6.1). We will now prove that $\mu^1 \in \mathfrak{C}$.

Since $\mu^1 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$ is a minimizer for (6.1), it follows that for any $\mu \in \mathcal{P}_2^\pi(\mathbb{R}^d)$,

$$(6.2) \quad \begin{aligned} \text{KL}(\mu^1|\pi) - \text{KL}(\mu|\pi) &\leq \frac{1}{\sigma} \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu^0, x)(\mu - \mu^1)(dx) \\ &\quad + \frac{1}{2\tau\sigma} \mathcal{W}_2(\mu, \mu^1) (\mathcal{W}_2(\mu, \mu^0) + \mathcal{W}_2(\mu^1, \mu^0)), \end{aligned}$$

By Theorem A.2, there exists a unique optimal coupling $\gamma^* \in \Gamma_o(\mu, \mu^1)$, and hence

$$(6.3) \quad \begin{aligned} \left| \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu^0, x)(\mu - \mu^1)(dx) \right| &= \left| \int_{\mathbb{R}^d \times \mathbb{R}^d} \left(\frac{\delta F}{\delta \mu}(\mu^0, x) - \frac{\delta F}{\delta \mu}(\mu^0, y) \right) \gamma^*(dx, dy) \right| \\ &\leq \int_{\mathbb{R}^d \times \mathbb{R}^d} \left| \frac{\delta F}{\delta \mu}(\mu^0, x) - \frac{\delta F}{\delta \mu}(\mu^0, y) \right| \gamma^*(dx, dy) \\ &\leq \left(\int_{\mathbb{R}^d \times \mathbb{R}^d} \left| \frac{\delta F}{\delta \mu}(\mu^0, x) - \frac{\delta F}{\delta \mu}(\mu^0, y) \right|^2 \gamma^*(dx, dy) \right)^{1/2} \\ &\leq L_F \left(\int_{\mathbb{R}^d \times \mathbb{R}^d} |x - y|^2 \gamma^*(dx, dy) \right)^{1/2} = L_F \mathcal{W}_2(\mu, \mu^1), \end{aligned}$$

where the penultimate inequality follows from the Cauchy-Schwarz inequality and the fact that $\gamma^* \in \mathcal{P}_2(\mathbb{R}^d \times \mathbb{R}^d)$, the last inequality follows from (2.2) in Assumption 2.2 and the last equality follows from optimality of γ^* . The rest of the proof is identical to the proof of Proposition 5.3. \square

From Proposition 6.3 it follows inductively that, for each $n \in \mathbb{N}$, given $\mu^n \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, the unique minimizer μ^{n+1} of (1.4) belongs to \mathfrak{C} . Hence, if $\mu^0 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, then $(\mu^n)_{n \in \mathbb{N}} \subset \mathfrak{C}$ along the scheme (1.4).

Since $\text{KL}(\cdot|\pi)$ admits a unique Wasserstein sub-differential and F is Wasserstein differentiable by Assumption 2.8, [1, Lemma 10.1.2] allows us to write the following first-order optimality condition for (1.4).

Lemma 6.4 (Optimality condition for (1.4)). *Let $F \in \mathcal{C}^1$, Assumption 2.2, 2.8 and (2.4) in Assumption 2.3 hold. Given $\mu^0 \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, then, for each $n \in \mathbb{N}$, the unique minimizer $\mu^{n+1} \in \mathfrak{C}$ of (1.4) satisfies*

$$\nabla_\mu F(\mu^n)(x) + \sigma \nabla \log \frac{d\mu^{n+1}}{d\pi}(x) = \frac{1}{\tau} \left(T_{\mu^{n+1}}^{\mu^n}(x) - x \right), \quad \text{for } \mu^{n+1}\text{-a.e. } x.$$

Proof. Identical to the proof of Lemma 5.4 once we replace F by $\int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu^n, x)(\mu - \mu^n)(dx)$ in J_n . \square

Since F is linearized in (1.4), convergence is not attainable without F being smooth. As in Euclidean geometry, where Lipschitz continuity of the gradients implies smoothness with respect to the squared Euclidean distance, an analogue implication holds in the Wasserstein space. In particular, as the following result shows, Assumption 2.9 implies that F is smooth with respect to \mathcal{W}_2^2 .

Lemma 6.5 (L_F -smoothness of F relative to \mathcal{W}_2^2). *Assume $F \in \mathcal{C}^1$ and Assumption 2.8, 2.9 hold. Then, for any $\mu', \mu \in \mathcal{P}_2(\mathbb{R}^d)$, it holds*

$$F(\mu') - F(\mu) - \left\langle \nabla_\mu F(\mu)(\cdot), P_\mu^{\mu'} - I_d \right\rangle_{L_\mu^2(\mathbb{R}^d)} \leq L'_F \left\| I_d - P_\mu^{\mu'} \right\|_{L_\mu^2(\mathbb{R}^d)}^2,$$

where $P_\mu^{\mu'} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the pushforward of μ onto μ' . If $P_\mu^{\mu'}$ is in fact an optimal transport map from μ to μ' , then $\|I_d - P_\mu^{\mu'}\|_{L_\mu^2(\mathbb{R}^d)}^2 = \mathcal{W}_2^2(\mu', \mu)$.

Proof. Let $\mu', \mu \in \mathcal{P}_2(\mathbb{R}^d)$. Then $\gamma = (I_d, P_\mu^{\mu'})_{\#} \mu$ is a coupling between μ and μ' . For any $\varepsilon \in [0, 1]$, set $\mu^\varepsilon = \mu + \varepsilon(\mu' - \mu)$. Then since $F \in \mathcal{C}^1$, it follows by Remark B.2 that

$$\begin{aligned} & F(\mu') - F(\mu) - \left\langle \nabla_\mu F(\mu)(\cdot), T_\mu^{\mu'} - I_d \right\rangle_{L_\mu^2(\mathbb{R}^d)} \\ &= \int_0^1 \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu^\varepsilon, x) (\mu' - \mu)(dx) d\varepsilon - \int_{\mathbb{R}^d} \nabla_\mu F(\mu)(x) \cdot (P_\mu^{\mu'}(x) - x) \mu(dx) \\ &= \int_0^1 \int_{\mathbb{R}^d \times \mathbb{R}^d} \left(\frac{\delta F}{\delta \mu}(\mu^\varepsilon, y) - \frac{\delta F}{\delta \mu}(\mu^\varepsilon, x) \right) \gamma(dx, dy) d\varepsilon - \int_{\mathbb{R}^d \times \mathbb{R}^d} \nabla_\mu F(\mu)(x) \cdot (y - x) \gamma(dx, dy). \end{aligned}$$

By Assumption 2.8, we have

$$\begin{aligned} & \int_0^1 \int_{\mathbb{R}^d \times \mathbb{R}^d} \left(\frac{\delta F}{\delta \mu}(\mu^\varepsilon, y) - \frac{\delta F}{\delta \mu}(\mu^\varepsilon, x) \right) \gamma^*(dx, dy) d\varepsilon \\ &= \int_0^1 \int_{\mathbb{R}^d \times \mathbb{R}^d} \int_0^1 \nabla_\mu F(\mu^\varepsilon)(x + \eta(y - x)) \cdot (y - x) d\eta \gamma(dx, dy) d\varepsilon. \end{aligned}$$

By Assumption 2.9, the Cauchy-Schwarz inequality and convexity of \mathcal{W}_2 , we obtain

$$\begin{aligned} & F(\mu') - F(\mu) - \left\langle \nabla_\mu F(\mu)(\cdot), P_\mu^{\mu'} - I_d \right\rangle_{L_\mu^2(\mathbb{R}^d)} \\ &= \int_0^1 \int_{\mathbb{R}^d \times \mathbb{R}^d} \int_0^1 (\nabla_\mu F(\mu^\varepsilon)(x + \eta(y - x)) - \nabla_\mu F(\mu)(x)) \cdot (y - x) d\eta \gamma(dx, dy) d\varepsilon \\ &\leq L'_F \int_0^1 \int_{\mathbb{R}^d \times \mathbb{R}^d} \int_0^1 (\eta|y - x| + \mathcal{W}_2(\mu^\varepsilon, \mu)) |y - x| d\eta \gamma(dx, dy) d\varepsilon \\ &\leq L'_F \int_0^1 \int_{\mathbb{R}^d \times \mathbb{R}^d} \int_0^1 (\eta|y - x| + \varepsilon \mathcal{W}_2(\mu', \mu)) |y - x| d\eta \gamma(dx, dy) d\varepsilon \\ &= \frac{L'_F}{2} \int_{\mathbb{R}^d \times \mathbb{R}^d} |y - x|^2 \gamma(dx, dy) + \frac{L'_F}{2} \mathcal{W}_2(\mu', \mu) \int_{\mathbb{R}^d \times \mathbb{R}^d} |y - x| \gamma(dx, dy) \\ &\leq L'_F \left\| I_d - P_\mu^{\mu'} \right\|_{L_\mu^2(\mathbb{R}^d)}^2. \end{aligned}$$

□

The following result is a consequence of the geodesic convexity of $\text{KL}(\cdot|\pi)$ in the Wasserstein space and combined with Lemma 6.5, it allows us to prove in Theorem 3.1 (ii) that $(F^\sigma(\mu^n))_n$ decreases along (1.4) as long as the step-size τ is small enough. It can also be viewed as a particular case of [31, Lemma 4]. However, our proof is slightly different.

Lemma 6.6. *Let Assumption 2.3 hold. For any $\mu', \mu \in \mathfrak{C}$, it holds*

$$\text{KL}(\mu'|\pi) - \text{KL}(\mu|\pi) - \left\langle \nabla \log \frac{d\mu'}{d\pi} \left(T_\mu^{\mu'} \right), T_\mu^{\mu'} - I_d \right\rangle_{L_\mu^2(\mathbb{R}^d)} \leq 0.$$

Proof. Let $\mu', \mu \in \mathfrak{C}$. Then since $\mathfrak{C} \subset \mathcal{P}_2^\pi(\mathbb{R}^d)$ and $\pi \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, it follows by Theorem A.2 that there exist unique μ -a.e. and μ' -a.e. optimal transport maps $T_\mu^{\mu'} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ from μ to μ' and $T_{\mu'}^\mu : \mathbb{R}^d \rightarrow \mathbb{R}^d$ from μ' to μ such that $T_{\mu'}^\mu \circ T_\mu^{\mu'} = I_d$, μ -a.e and $T_\mu^{\mu'} \circ T_{\mu'}^\mu = I_d$, μ' -a.e..

Now we show that, for any $\varepsilon \in (0, 1)$, $I_d + \varepsilon (T_{\mu'}^\mu - I_d)$ is the unique optimal transport map from μ' to $(I_d + \varepsilon (T_{\mu'}^\mu - I_d))_{\#} \mu'$. First, we observe that

$$\int_{\mathbb{R}^d} |x + \varepsilon (T_{\mu'}^\mu(x) - x)|^2 \mu'(dx) \leq 2 \int_{\mathbb{R}^d} |x|^2 \mu'(dx) + 2\varepsilon^2 \mathcal{W}_2^2(\mu', \mu) < \infty.$$

Hence, $(I_d + \varepsilon (T_{\mu'}^\mu - I_d))_{\#} \mu' \in \mathcal{P}_2(\mathbb{R}^d)$. Therefore, by Theorem A.4, it suffices to show that $I_d + \varepsilon (T_{\mu'}^\mu - I_d)$ is the gradient of a convex differentiable μ' -a.e. function. By Theorem A.2, we have that $T_{\mu'}^\mu(x) = \nabla \varphi(x)$ μ' -a.e. for a convex function $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}$. Hence, for μ' -a.e. x ,

$$x + \varepsilon (T_{\mu'}^\mu(x) - x) = (1 - \varepsilon)x + \varepsilon T_{\mu'}^\mu(x) = \nabla \left((1 - \varepsilon) \frac{|x|^2}{2} + \varepsilon \varphi(x) \right),$$

where the map $\mathbb{R}^d \ni x \mapsto (1 - \varepsilon) \frac{|x|^2}{2} + \varepsilon \varphi(x) \in \mathbb{R}$ is convex and μ' -a.e. differentiable for all $\varepsilon \in (0, 1)$. The uniqueness of $I_d + \varepsilon (T_{\mu'}^\mu - I_d)$ follows from Theorem A.2 since $\mu' \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$ and the fact that $(I_d + \varepsilon (T_{\mu'}^\mu - I_d))_{\#} \mu' \in \mathcal{P}_2(\mathbb{R}^d)$.

By Theorem A.5, since $\mu' \in \mathfrak{C}$,

$$\partial^- \text{KL}(\mu'|\pi) = \left\{ \nabla \log \frac{d\mu'}{d\pi} \right\}.$$

Hence, by (B.2), we have

$$(6.4) \quad \begin{aligned} \text{KL} \left((I_d + \varepsilon (T_{\mu'}^\mu - I_d))_{\#} \mu' \middle| \pi \right) &\geq \text{KL}(\mu'|\pi) + \varepsilon \int_{\mathbb{R}^d} \nabla \log \frac{d\mu'}{d\pi}(x) \cdot (T_{\mu'}^\mu(x) - x) \mu'(dx) \\ &\quad + o \left(\mathcal{W}_2 \left(\mu', (I_d + \varepsilon (T_{\mu'}^\mu - I_d))_{\#} \mu' \right) \right). \end{aligned}$$

Now, by Corollary A.3, we have

$$(6.5) \quad \mathcal{W}_2 \left(\mu', (I_d + \varepsilon (T_{\mu'}^\mu - I_d))_{\#} \mu' \right) = \varepsilon \mathcal{W}_2(\mu', \mu).$$

By (2.5) in Assumption 2.3, π is log-concave, thus by [1, Theorem 9.4.10] the relative entropy $\text{KL}(\cdot|\pi)$ is geodesically convex, and hence

$$(6.6) \quad \text{KL} \left((I_d + \varepsilon (T_{\mu'}^\mu - I_d))_{\#} \mu' \middle| \pi \right) \leq (1 - \varepsilon) \text{KL}(\mu'|\pi) + \varepsilon \text{KL}(\mu|\pi).$$

Combining (6.4), (6.5) and (6.6) and using (A.1) gives

$$\begin{aligned} \text{KL}(\mu|\pi) - \text{KL}(\mu'|\pi) &\geq \int_{\mathbb{R}^d} \nabla \log \frac{d\mu'}{d\pi}(x) \cdot (T_{\mu'}^\mu(x) - x) \mu'(dx) + \frac{o(\varepsilon)}{\varepsilon} \\ &= \int_{\mathbb{R}^d} \nabla \log \frac{d\mu'}{d\pi}(x) \cdot (T_{\mu'}^\mu(x) - x) \left(T_{\mu'}^\mu \right)_{\#} \mu(dx) + \frac{o(\varepsilon)}{\varepsilon} \\ &= \int_{\mathbb{R}^d} \nabla \log \frac{d\mu'}{d\pi} \left(T_{\mu'}^\mu(x) \right) \cdot \left(x - T_{\mu'}^\mu(x) \right) \mu(dx) + \frac{o(\varepsilon)}{\varepsilon}. \end{aligned}$$

Sending $\varepsilon \rightarrow 0$ and rearranging gives the conclusion. \square

7. PROXIMAL GRADIENT SCHEME

In this section, we present the auxilliary results needed for the proof of (iii) in Theorem 3.1. Before applying the same proof of Theorem 5.1 and 6.1 to show existence and uniqueness of a minimizer in the JKO step of (1.5), we prove that the pushforward of μ^n by $I_d - \tau \nabla_\mu F(\mu^n)(\cdot)$ is an optimal transport from μ^n to $\nu^{n+1} = (I_d - \tau \nabla_\mu F(\mu^n)(\cdot))_{\#} \mu^n$, and moreover that $\nu^{n+1} \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$. As we saw in the proofs of Theorem 5.1 and 6.1, the previous step in the JKO update, in this case ν^{n+1} , needs to be absolutely continuous. This is proved in the following lemma, which is a generalization of [31, Lemma 2].

Lemma 7.1. *Let Assumption 2.8, 2.9 hold. Let $\mu \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, $\sigma > 0$ and $\nu = (I_d - \tau \nabla_\mu F(\mu)(\cdot))_{\#} \mu$. If $\tau < \frac{1}{L'_F}$, the optimal transport map from μ to ν is given by*

$$T_\mu^\nu = I_d - \tau \nabla_\mu F(\mu)(\cdot).$$

Moreover, $\nu \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$.

Proof. Note that Assumption 2.8 together with that fact that $\mu \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$ imply that

$$\int_{\mathbb{R}^d} |x - \tau \nabla_\mu F(\mu)(x)|^2 \mu(dx) < \infty,$$

and hence $\nu \in \mathcal{P}_2(\mathbb{R}^d)$.

Since the map $I_d - \tau \nabla_\mu F(\mu)(\cdot)$ is a pushforward from μ to ν , by Theorem A.4, it suffices to show that $I_d - \tau \nabla_\mu F(\mu)(\cdot)$ can be written as the gradient of a convex function.

Let $u(x) := \frac{1}{2}|x|^2 - \tau \frac{\delta F}{\delta \mu}(\mu, x)$. Then, for any $x \in \mathbb{R}^d$, $\nabla u(x) = x - \tau \nabla_\mu F(\mu)(x)$. Moreover,

$$\begin{aligned} (x - y) \cdot (\nabla u(x) - \nabla u(y)) &= (x - y) \cdot (x - y - \tau (\nabla_\mu F(\mu)(x) - \nabla_\mu F(\mu)(y))) \\ &= |x - y|^2 - \tau (x - y) \cdot (\nabla_\mu F(\mu)(x) - \nabla_\mu F(\mu)(y)) \\ &\geq |x - y|^2 - \tau |x - y| |\nabla_\mu F(\mu)(x) - \nabla_\mu F(\mu)(y)| \\ &\geq (1 - \tau L'_F) |x - y|^2. \end{aligned}$$

Since by assumption $\tau < \frac{1}{L'_F}$, it follows that u is $(1 - \tau L'_F)$ -strongly convex and moreover ∇u is injective. By strong convexity of u and Theorem A.4, we obtain that the pushforward from μ to ν via ∇u is an optimal transport map, and we denote it by

$$T_\mu^\nu = I_d - \tau \nabla_\mu F(\mu)(\cdot).$$

By injectivity of ∇u and strong convexity of u , we obtain from [1, Lemma 5.5.3] that $\nu \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$. \square

Now, since $\nu^{n+1} \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$ for a sufficiently small step-size τ , the existence and uniqueness of a minimizer for (1.5) is a consequence of [20, Proposition 4.1].

Theorem 7.2 (Existence and uniqueness of minimizer for (1.5)). *Let Assumption 2.8, 2.9 and (2.4) in Assumption 2.3 hold. If $\tau < \frac{1}{L'_F}$, given $\mu^0 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, there exists a unique minimizer $\mu^1 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$ of*

$$(7.1) \quad \mathcal{P}_2^\lambda(\mathbb{R}^d) \ni \mu \mapsto \mathcal{H}(\mu) := \sigma \text{KL}(\mu|\pi) + \frac{1}{2\tau} \mathcal{W}_2^2(\mu, \nu^1).$$

Proof. From Lemma 7.1 it follows that given $\mu^0 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, we obtain $\nu^1 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$. Hence, by [20, Proposition 4.1], there exists a unique minimizer $\mu^1 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$ for \mathcal{H} . \square

From Lemma 7.1 and Theorem 7.2, it follows inductively that, for each $n \in \mathbb{N}$, given $\mu^n \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, we obtain $\nu^{n+1} \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, and hence $\mu^{n+1} \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$. Therefore, if $\mu^0 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, then $(\nu^n, \mu^n)_{n \in \mathbb{N}} \subset \mathcal{P}_2^\lambda(\mathbb{R}^d) \times \mathcal{P}_2^\lambda(\mathbb{R}^d)$ along the scheme (1.5). Therefore, via Theorem A.2, we obtain

Corollary 7.3 (Existence of optimal transport maps along (1.5)). *Let Assumption 2.8, 2.9 and (2.4) in Assumption 2.3 hold. If $\tau < \frac{1}{L_F}$, given $\mu^0 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, there exist unique μ^n -a.e. and ν^{n+1} -a.e. optimal transport maps $T_{\mu^n}^{\nu^{n+1}} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $T_{\nu^{n+1}}^{\mu^n} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ from μ^n to ν^{n+1} and from ν^{n+1} to μ^n , respectively, given by $T_{\mu^n}^{\nu^{n+1}} = I_d - \tau \nabla_\mu F(\mu^n)(\cdot)$ and $T_{\nu^{n+1}}^{\mu^n} = (I_d - \tau \nabla_\mu F(\mu^n)(\cdot))^{-1}$. Moreover, there also exist unique ν^{n+1} -a.e. and μ^{n+1} -a.e. optimal transport maps $T_{\nu^{n+1}}^{\mu^{n+1}} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $T_{\mu^{n+1}}^{\nu^{n+1}} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that $T_{\nu^{n+1}}^{\mu^{n+1}} \circ T_{\mu^{n+1}}^{\nu^{n+1}} = I_d$, μ^{n+1} -a.e. and $T_{\mu^{n+1}}^{\nu^{n+1}} \circ T_{\nu^{n+1}}^{\mu^{n+1}} = I_d$, ν^{n+1} -a.e..*

The following proposition is a particular case of Proposition 5.3 and 6.3 and guarantees that H admits a unique Wasserstein sub-differential given by $\nabla \log \frac{d\mu^n}{d\lambda}$, where μ^n is the iterate generated by (1.5) at each step $n \geq 1$.

Proposition 7.4 (Wasserstein sub-differentiability class for $\text{KL}(\cdot|\pi)$ along (1.5)). *Let Assumption 2.8, 2.9 and (2.4) in Assumption 2.3 hold. If $\tau < \frac{1}{L_F}$, given $\mu^0 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, the unique minimizer $\mu^1 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$ of (7.1) belongs to \mathfrak{C} .*

Proof. Since $\mu^0 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, Theorem 7.2 guarantees the existence and uniqueness of a minimizer $\mu^1 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$ for (7.1). Following either the proof of Proposition 5.3 or Proposition 6.3 will show that $\mu^1 \in \mathfrak{C}$. \square

Lemma 7.5 (Optimality condition for (1.5)). *Let Assumption 2.8, 2.9 and (2.4) in Assumption 2.3 hold. If $\tau < \frac{1}{L_F}$, given $\mu^0 \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, then, for each $n \in \mathbb{N}$, the unique minimizer $\mu^{n+1} \in \mathfrak{C}$ of the minimization step in (1.5) satisfies*

$$\sigma \nabla \log \frac{d\mu^{n+1}}{d\pi}(x) = \frac{1}{\tau} \left(T_{\mu^{n+1}}^{\nu^{n+1}}(x) - x \right), \quad \text{for } \mu^{n+1}\text{-a.e. } x.$$

Proof. The result follows from either Lemma 5.4 or Lemma 6.4 with $F = 0$. \square

APPENDIX A. OPTIMAL TRANSPORT

In this appendix, we recall the fundamental results from optimal transport that are used throughout the paper.

Definition A.1 (Pushforward of a measure by a map). Let $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be a $\mathcal{B}(\mathbb{R}^d)$ -measurable map. Then, for every $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, we denote by $T_{\#}\mu \in \mathcal{P}_2(\mathbb{R}^d)$ the pushforward measure of μ by T , characterized by

$$(A.1) \quad \int_{\mathbb{R}^d} f(T(x))\mu(dx) = \int_{\mathbb{R}^d} f(y) (T_{\#}\mu)(dy), \quad \text{for any measurable bounded function } f.$$

Consider the 2-Wasserstein distance $\mathcal{W}_2 : \mathcal{P}_2(\mathbb{R}^d) \times \mathcal{P}_2(\mathbb{R}^d) \rightarrow [0, \infty)$, defined by

$$(A.2) \quad \mathcal{W}_2(\mu, \nu) := \left(\inf_{\gamma \in \Gamma(\mu, \nu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} |x - y|^2 \gamma(dx, dy) \right)^{\frac{1}{2}},$$

where $\Gamma(\mu, \nu) := \{\gamma \in \mathcal{P}_2(\mathbb{R}^d \times \mathbb{R}^d) : (P_x)_{\#}\gamma = \mu, (P_y)_{\#}\gamma = \nu\}$ is the set of couplings between μ and ν , where $P_x : (x, y) \mapsto x$ and $P_y : (x, y) \mapsto y$ are the projections onto the first

and second component, respectively. The set of optimal couplings for which the infimum is attained in (A.2) is denoted by $\Gamma_o(\mu, \nu) := \left\{ \bar{\gamma} \in \Gamma(\mu, \nu) : \mathcal{W}_2^2(\mu, \nu) = \int_{\mathbb{R}^d \times \mathbb{R}^d} |x - y|^2 \bar{\gamma}(dx, dy) \right\}$.

Now we recall a standard result from optimal transport (see e.g. [15, Theorem 4.5] and also [1, 33]), which shall be essential throughout the paper.

Theorem A.2. *Let $\mu \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$ and $\nu \in \mathcal{P}_2(\mathbb{R}^d)$. Then*

- (i) *there exists a unique optimal coupling $\gamma^* := (I_d, T_\mu^\nu)_\# \mu$ which minimizes (A.2), where $T_\mu^\nu : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the unique μ -almost everywhere (a.e.) optimal transport map from μ to ν .*
- (ii) *Moreover, $T_\mu^\nu(x) = \nabla \varphi(x)$ μ -a.e. for a convex function $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}$ with $\varphi(x) := \frac{1}{2}|x|^2 - \psi(x)$, where ψ is a c-concave function.²*
- (iii) *If $\nu \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$, then $\gamma^* = (T_\nu^\mu, I_d)_\# \nu$, where $T_\nu^\mu : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the unique ν -a.e. optimal transport map such that*

$$T_\nu^\mu \circ T_\mu^\nu = I_d, \quad \mu\text{-a.e. and } T_\mu^\nu \circ T_\nu^\mu = I_d, \quad \nu\text{-a.e.}$$

As a consequence of Theorem A.2, we have the following

Corollary A.3. *Let $\mu, \nu \in \mathcal{P}_2^\lambda(\mathbb{R}^d)$. Then*

$$\mathcal{W}_2^2(\mu, \nu) = \int_{\mathbb{R}^d} |x - T_\mu^\nu(x)|^2 \mu(dx) = \int_{\mathbb{R}^d} |x - T_\nu^\mu(x)|^2 \nu(dx).$$

Theorem A.4 ([32, Theorem 1.48]). *Suppose $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ and that $u : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex and differentiable μ -a.e. Set $T := \nabla u$ and suppose $\int_{\mathbb{R}^d} |T(x)|^2 \mu(dx) < \infty$. Then T is an optimal transport map from μ to $T_\# \mu$.*

Theorem A.5 (Subdifferential of $\text{KL}(\cdot|\pi)$; [1, Theorem 10.4.9]). *The relative entropy $\text{KL}(\cdot|\pi)$ has finite slope at $\mu \in \mathcal{P}_2^\pi(\mathbb{R}^d)$, i.e.,*

$$|\partial \text{KL}(\cdot|\pi)|(\mu) := \limsup_{\nu \rightarrow \mu} \frac{(\text{KL}(\mu|\pi) - \text{KL}(\nu|\pi))_+}{\mathcal{W}_2(\nu, \mu)} < \infty$$

if and only if $\frac{d\mu}{d\pi} \in W_{\lambda, \text{loc}}^{1,1}(\mathbb{R}^d)$ and $\nabla \log \frac{d\mu}{d\pi} \in L_\mu^2(\mathbb{R}^d)$. In this case, $I(\mu|\pi) = |\partial \text{KL}(\cdot|\pi)|^2(\mu)$ and $\partial^- \text{KL}(\mu|\pi) = \left\{ \nabla \log \frac{d\mu}{d\pi} \right\}$ (cf. (B.2) in Definition B.3).

APPENDIX B. DIFFERENTIAL CALCULUS ON $\mathcal{P}_2(\mathbb{R}^d)$

In this appendix, we recall the notions of linear functional (flat) differentiability [6] and Wasserstein differentiability [7] used throughout the paper.

Definition B.1 (Flat differentiability on $\mathcal{P}_2(\mathbb{R}^d)$). We say a function $F : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ is in \mathcal{C}^1 , if there exists a continuous function $\frac{\delta F}{\delta \mu} : \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d \rightarrow \mathbb{R}$, with respect to the product topology on $\mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d$, called the flat derivative of F , for which there exists $\kappa > 0$ such that for all $(\mu, x) \in \mathcal{P}_2(\mathbb{R}^d) \times \mathbb{R}^d$, $\left| \frac{\delta F}{\delta \mu}(\mu, x) \right| \leq \kappa(1 + |x|^2)$, and for all $\mu' \in \mathcal{P}_2(\mathbb{R}^d)$,

$$(B.1) \quad \lim_{\varepsilon \searrow 0} \frac{F(\mu^\varepsilon) - F(\mu)}{\varepsilon} = \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu, x)(\mu' - \mu)(dx), \quad \text{with } \mu^\varepsilon = \mu + \varepsilon(\mu' - \mu),$$

and $\int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu, x) \mu(dx) = 0$.

²A function $\xi : \mathbb{R}^d \rightarrow \mathbb{R}$ is c-concave if and only if $x \mapsto \frac{1}{2}|x|^2 - \xi(x)$ is convex and lower semi-continuous.

Remark B.2. One can show that if $F : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}^d$ admits a flat derivative $\frac{\delta F}{\delta \mu}$, then for all $\mu, \mu' \in \mathcal{P}_2(\mathbb{R}^d)$, the function $[0, 1] \ni \varepsilon \mapsto f(\mu^\varepsilon)$ is continuous on $[0, 1]$ and differentiable on $(0, 1)$ with derivative $\frac{d}{d\varepsilon} f(\mu^\varepsilon) = \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu^\varepsilon, x)(\mu' - \mu)(dx)$ (see [21, Theorem 2.3]). Hence, by the fundamental theorem of calculus, $F(\mu') - F(\mu) = \int_0^1 \int_{\mathbb{R}^d} \frac{\delta F}{\delta \mu}(\mu^\varepsilon, x)(\mu' - \mu)(dx)d\varepsilon$, provided that $\varepsilon \mapsto \int \frac{\delta F}{\delta \mu}(\mu^\varepsilon, x)(\mu' - \mu)(dx)$ is integrable.

Recall that the tangent space of $\mathcal{P}_2(\mathbb{R}^d)$ at $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ is defined as

$$\mathcal{T}_\mu \mathcal{P}_2(\mathbb{R}^d) = \overline{\{\nabla \psi : \psi \in C_c^\infty(\mathbb{R}^d)\}} \subset L_\mu^2(\mathbb{R}^d),$$

where the closure is taken in $L_\mu^2(\mathbb{R}^d)$, see [1, Definition 8.4.1], and $C_c^\infty(\mathbb{R}^d)$ denotes the space of smooth functions with compact support in \mathbb{R}^d .

Definition B.3 (Wasserstein sub- and super-differential on $\mathcal{P}_2(\mathbb{R}^d)$). Let $F : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ and let $\mu \in \mathcal{P}_2(\mathbb{R}^d)$. Then

- (i) a map $\xi \in \mathcal{T}_\mu \mathcal{P}_2(\mathbb{R}^d)$ belongs to the sub-differential $\partial^- F(\mu)$ of F at $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ if for all $\mu' \in \mathcal{P}_2(\mathbb{R}^d)$,

$$(B.2) \quad F(\mu') \geq F(\mu) + \sup_{\gamma \in \Gamma_o(\mu, \mu')} \int_{\mathbb{R}^d \times \mathbb{R}^d} \langle \xi(x), y - x \rangle d\gamma(x, y) + o(\mathcal{W}_2(\mu, \mu')).$$

If $\partial^- F(\mu) \neq \emptyset$, we say the function F is Wasserstein sub-differentiable at μ .

- (ii) A map $\xi \in \mathcal{T}_\mu \mathcal{P}_2(\mathbb{R}^d)$ belongs to the super-differential $\partial^+ F(\mu)$ of F at $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ if $-\xi \in \partial^-(-F)(\mu)$. If $\partial^+ F(\mu) \neq \emptyset$, we say the function F is Wasserstein super-differentiable at μ .

Then, we say that a function is Wasserstein differentiable if it admits sub- and super-differentials which coincide.

Definition B.4 (Wasserstein differentiability on $\mathcal{P}_2(\mathbb{R}^d)$). We say that a function $F : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ is Wasserstein differentiable at $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ if $\partial^- F(\mu) \cap \partial^+ F(\mu) \neq \emptyset$.

If $F : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ is Wasserstein differentiable at $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ (cf. Definition B.4), then by [7, Proposition 5.63], there exists a unique map $\nabla_\mu F(\mu) \in \mathcal{T}_\mu \mathcal{P}_2(\mathbb{R}^d)$ such that $\partial^- F(\mu) = \partial^+ F(\mu) = \{\nabla_\mu F(\mu)\}$, called the Wasserstein gradient of F at $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, satisfying for any $\mu' \in \mathcal{P}_2(\mathbb{R}^d)$, and $\gamma \in \Gamma_o(\mu, \mu')$,

$$F(\mu') = F(\mu) + \int_{\mathbb{R}^d \times \mathbb{R}^d} \langle \nabla_\mu F(\mu)(x), y - x \rangle d\gamma(x, y) + o(\mathcal{W}_2(\mu, \mu')).$$

ACKNOWLEDGEMENTS

R-AL was supported by the EPSRC Centre for Doctoral Training in Mathematical Modelling, Analysis and Computation (MAC-MIGS) funded by the UK Engineering and Physical Sciences Research Council (grant EP/S023291/1), Heriot-Watt University and the University of Edinburgh. LS acknowledges the support of the UKRI Prosperity Partnership Scheme (FAIR) under EPSRC Grant EP/V056883/1 and the Alan Turing Institute.

REFERENCES

- [1] L. Ambrosio, N. Gigli, and G. Savare. *Gradient Flows: In Metric Spaces and in the Space of Probability Measures*. Lectures in Mathematics. ETH Zürich. Birkhäuser Basel, 2008.
- [2] M. Arbel, A. Korba, A. Salim, and A. Gretton. Maximum mean discrepancy gradient flow. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

- [3] D. Bakry and M. Émery. Diffusions hypercontractives. *Séminaire de probabilités de Strasbourg*, 19:177–206, 1985.
- [4] J.-D. Benamou, G. Carlier, Q. Mérigot, and É. Oudet. Discretization of functionals involving the Monge–Ampère operator. *Numerische Mathematik*, 134:611–636, 2014.
- [5] V. Bogachev. *Measure Theory: Volume 1*. Springer Berlin Heidelberg, 2007.
- [6] P. Cardaliaguet, F. Delarue, J. Lasry, and P. Lions. *The Master Equation and the Convergence Problem in Mean Field Games*. Annals of Mathematics Studies. Princeton University Press, 2019.
- [7] R. A. Carmona and F. Delarue. *Probabilistic Theory of Mean Field Games with Applications I: Mean Field FBSDEs, Control, and Games*. Springer International Publishing, 2018.
- [8] P. Cattiaux and A. Guillin. Functional inequalities for perturbed measures with applications to log-concave measures and to some Bayesian problems. *Bernoulli*, 28(4):2294 – 2321, 2022.
- [9] F. Chen, Z. Ren, and S. Wang. Entropic fictitious play for mean field optimization problem. *Journal of Machine Learning Research*, 24(211):1–36, 2023.
- [10] F. Chen, Z. Ren, and S. Wang. Uniform-in-time propagation of chaos for mean field Langevin dynamics, 2023. arXiv:2212.03050.
- [11] S. Chewi, A. Nitanda, and M. S. Zhang. Uniform-in- n log-Sobolev inequality for the mean-field Langevin dynamics with convex energy, 2024. arXiv:2409.10440.
- [12] L. Chizat. Mean-field Langevin dynamics : Exponential convergence and annealing. *Transactions on Machine Learning Research*, 2022.
- [13] L. Chizat and F. R. Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. In *NeurIPS*, 2018.
- [14] D. Drusvyatskiy and C. Paquette. Efficiency of minimizing compositions of convex functions and smooth maps. *Mathematical Programming*, pages 1–56, 2016.
- [15] W. Gangbo and R. J. McCann. The geometry of optimal transportation. *Acta Mathematica*, 177(2):113 – 161, 1996.
- [16] R. Holley and D. W. Stroock. Logarithmic Sobolev inequalities and stochastic Ising models. *Journal of Statistical Physics*, 46:1159–1194, 1987.
- [17] K. Hu, Z. Ren, D. Šiška, and L. Szpruch. Mean-field Langevin dynamics and energy landscape of neural networks. *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, 57(4):2043 – 2065, 2021.
- [18] Y.-J. Huang and Z. Malik. Generative modeling by minimizing the Wasserstein-2 loss, 2024. arXiv:2406.13619.
- [19] M. L. Jean-David Benamou, Guillaume Carlier. An augmented Lagrangian approach to Wasserstein gradient flows and applications. *ESAIM: Proceedings and Surveys*, 54(1):1–17, 2016.
- [20] R. Jordan, D. Kinderlehrer, and F. Otto. The variational formulation of the Fokker–Planck equation. *SIAM Journal on Mathematical Analysis*, 29(1):1–17, 1998.
- [21] B. Jourdain and A. Tse. Central limit theorem over non-linear functionals of empirical measures with applications to the mean-field fluctuation of interacting diffusions. *Electronic Journal of Probability*, 26:1 – 34, 2021.
- [22] J.-M. Leahy, B. Kerimkulov, D. Šiška, and L. Szpruch. Convergence of policy gradient for entropy regularized MDPs with neural network approximation in the mean-field regime. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 12222–12252. PMLR, 17–23 Jul 2022.
- [23] H. P. H. Luu, H. Yu, B. Williams, P. Mikkola, M. Hartmann, K. Puolamäki, and A. Klami. Non-geodesically-convex optimization in the Wasserstein space. In *The 38th Annual Conference on Neural Information Processing Systems*, 2024.
- [24] S. Mei, A. Montanari, and P.-M. Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences of the United States of America*, 115:E7665 – E7671, 2018.
- [25] P. Monmarché, Z. Ren, and S. Wang. Time-uniform log-Sobolev inequalities and applications to propagation of chaos. *Electronic Journal of Probability*, 29(none):1 – 38, 2024.
- [26] A. Nitanda and T. Suzuki. Stochastic particle gradient descent for infinite ensembles, 2017. arXiv:1712.05438.
- [27] A. Nitanda, D. Wu, and T. Suzuki. Convex analysis of the mean field Langevin dynamics. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, 2022.
- [28] F. Otto and C. Villani. Generalization of an inequality by Talagrand and links with the logarithmic Sobolev inequality. *Journal of Functional Analysis*, 173(2):361–400, 2000.

- [29] N. Parikh and S. Boyd. Proximal algorithms. *Found. Trends Optim.*, 1(3):127–239, Jan. 2014.
- [30] G. M. Rotskoff and E. Vanden-Eijnden. Trainability and accuracy of artificial neural networks: An interacting particle system approach. *Communications on Pure and Applied Mathematics*, 75, 2018.
- [31] A. Salim, A. Korba, and G. Luise. The Wasserstein proximal gradient algorithm. In *Advances in Neural Information Processing Systems*, volume 33, pages 12356–12366. Curran Associates, Inc., 2020.
- [32] F. Santambrogio. *Optimal Transport for Applied Mathematicians: Calculus of Variations, PDEs, and Modeling*. Progress in Nonlinear Differential Equations and Their Applications. Springer International Publishing, 2015.
- [33] F. Santambrogio. Euclidean, Metric, and Wasserstein gradient flows: an overview. *Bulletin of Mathematical Sciences*, 7, 09 2016.
- [34] J. Sirignano and K. Spiliopoulos. Mean field analysis of neural networks: A law of large numbers. *SIAM Journal on Applied Mathematics*, 80(2):725–752, 2020.
- [35] A. Teter, I. Nodouzi, and A. Halder. Proximal mean field learning in shallow neural networks. *Transactions on Machine Learning Research*, 2024.
- [36] C. Villani. *Optimal Transport: Old and New*. Grundlehren der mathematischen Wissenschaften. Springer Berlin Heidelberg, 2008.
- [37] S. Wang. Uniform log-Sobolev inequalities for mean field particles with flat-convex energy, 2024. arXiv:2408.03283.
- [38] A. Wibisono. Sampling as optimization in the space of measures: The Langevin dynamics as a composite optimization problem. In *Conference on Learning Theory*, 2018.
- [39] Y. Xu and Q. Li. Forward-Euler time-discretization for Wasserstein gradient flows can be wrong, 2024. arXiv:2406.08209.
- [40] R. Zhang, C. Chen, C. Li, and L. Carin. Policy optimization as Wasserstein gradient flows. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5737–5746. PMLR, 10–15 Jul 2018.

SCHOOL OF MATHEMATICAL AND COMPUTER SCIENCES, HERIOT-WATT UNIVERSITY, EDINBURGH, UK, AND MAXWELL INSTITUTE FOR MATHEMATICAL SCIENCES, EDINBURGH, UK

Email address: r12029@hw.ac.uk

SCHOOL OF MATHEMATICAL AND COMPUTER SCIENCES, HERIOT-WATT UNIVERSITY, EDINBURGH, UK, AND MAXWELL INSTITUTE FOR MATHEMATICAL SCIENCES, EDINBURGH, UK

Email address: m.majka@hw.ac.uk

SCHOOL OF MATHEMATICS, UNIVERSITY OF EDINBURGH, UK, AND SIMTOPIA, UK

Email address: d.siska@ed.ac.uk

SCHOOL OF MATHEMATICS, UNIVERSITY OF EDINBURGH, UK, THE ALAN TURING INSTITUTE, UK AND SIMTOPIA, UK

Email address: l.szpruch@ed.ac.uk