PROVABLY CONVERGENT STOCHASTIC FIXED-POINT ALGORITHM FOR FREE-SUPPORT WASSERSTEIN BARYCENTER OF CONTINUOUS NON-PARAMETRIC MEASURES

ZEYI CHEN, ARIEL NEUFELD, AND QIKUN XIANG

ABSTRACT. We propose a provably convergent algorithm for approximating the 2-Wasserstein barycenter of continuous non-parametric probability measures. Our algorithm is inspired by the fixed-point iterative scheme of Álvarez-Esteban et al. (2016) whose convergence to the 2-Wasserstein barycenter relies on obtaining exact optimal transport (OT) maps. However, typically in practice, OT maps are only approximately computed and exact computation of OT maps between continuous probability measures is only tractable for certain restrictive parametric families. To circumvent the need to compute exact OT maps between general non-parametric measures, we develop a tailored iterative scheme that utilizes consistent estimators of the OT maps instead of the exact OT maps. This gives rise to a computationally tractable stochastic fixed-point algorithm which is provably convergent to the 2-Wasserstein barycenter. Our algorithm remarkably does not restrict the support of the 2-Wasserstein barycenter to be any fixed finite set and can be implemented in a distributed computing environment, which makes it suitable for large-scale data aggregation problems. In our numerical experiments, we propose a method of generating non-trivial instances of 2-Wasserstein barycenter problems where the ground-truth barycenter measure is known. Our numerical results showcase the capability of our algorithm in developing high-quality approximations of the 2-Wasserstein barycenter, as well as its superiority over state-of-the-art methods based on generative neural networks in terms of accuracy, stability, and efficiency.

Keywords: Wasserstein barycenter, optimal transport, information aggregation, transportation map estimation

1. INTRODUCTION

Aggregating information from multiple heterogeneous data sources is broadly encountered in many application scenarios. Typical instances include forming group consensus from expert judgements or forecasts in decision analysis [50, 54], combining subsample posteriors for large datasets in Bayesian inference [6, 44], pooling dependent samples under data scarcity in data-driven optimization [33, 57], adapting informative signals across domains in transfer learning [7, 64], etc. A common technique in information aggregation is to summarize the data characteristics by utilizing some type of barycenter (also known as the Karcher–Fréchet mean [28, 36]) which retrieves a formal notion of "centroid" of points in certain metric space. In this paper, we are interested in the particular case of information aggregation where data from K > 2 sources are represented by probability measures ν_1, \ldots, ν_K on \mathbb{R}^d with $d \in \mathbb{N}$, and are aggregated via their 2-Wasserstein barycenter (W_2 -barycenter) [1] defined as follows.

Definition 1.1 (W_2 -distance and W_2 -barycenter [1]). The 2-Wasserstein distance, or W_2 -distance, between two probability measures $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$ with finite second moments is defined via the following optimal transport problem (see, e.g., [69]) with squared-distance cost:

$$\mathcal{W}_{2}(\mu,\nu) := \left(\inf_{\pi \in \Pi(\mu,\nu)} \int_{\mathbb{R}^{d} \times \mathbb{R}^{d}} \|\boldsymbol{x} - \boldsymbol{y}\|^{2} \, \pi(\mathrm{d}\boldsymbol{x},\mathrm{d}\boldsymbol{y})\right)^{\frac{1}{2}},\tag{1.1}$$

where $\Pi(\mu, \nu)$ denotes the set of couplings between μ and ν (see Definition 2.2). For $\nu_1, \ldots, \nu_K \in \mathcal{P}_2(\mathbb{R}^d)$, weights $w_1 > 0, \ldots, w_K > 0$ satisfying $\sum_{k=1}^{K} w_k = 1$, and for any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, let $V(\mu)$ denote the convex combination of the squared \mathcal{W}_2 -distances between μ and ν_1, \ldots, ν_K given by

$$V(\mu) := \sum_{k=1}^{K} w_k \mathcal{W}_2(\mu, \nu_k)^2.$$
(1.2)

Then, $\bar{\mu} \in \mathcal{P}_2(\mathbb{R}^d)$ is called a \mathcal{W}_2 -barycenter of ν_1, \ldots, ν_K with weights w_1, \ldots, w_K if

$$\bar{\mu} \in \underset{\mu \in \mathcal{P}_2(\mathbb{R}^d)}{\operatorname{arg\,min}} V(\mu).$$

Date: June 2, 2025.

In words, the W_2 -distance between two probability measures is defined as the minimal transportation cost of moving probability mass from one to the other under the squared-distance cost function. This induces a metric on the space of probability measures with finite second moments that metrizes the weak convergence; see, e.g., [69, Theorem 6.9]. Due to the appealing geometric intuition and statistical properties of the W_2 -distance, the W_2 -barycenter has been serving as a powerful tool in widespread applications in terms of distribution aggregation and representation tasks, including but not limited to computer graphics [53], machine learning [23, 45], theoretical economics [48, 50], Bayesian statistics [60, 61], network analysis [59], etc. However, it is well-known that the computation of the Wasserstein barycenter for general non-parametric continuous measures suffers from poor scalability, which has become the major bottleneck for its practical usage. In fact, even in the restrictive case of aggregating discrete measures, it has been proved by Altschuler and Boix-Adserà [2] that the time complexity for computing the W_2 -barycenter grows exponentially with the dimension, and thus the problem is NP-hard. A common strategy to approximate the W_2 -barycenter is to parametrize it via a discrete measure supported on fixed atoms, which transforms the problem into optimizing the histogram weights over a finite-dimensional probability simplex; see, e.g., [51, Chapter 6] and the references therein. Nevertheless, all such fixed-support algorithms have poor scalability in high dimensions due to prohibitive computational burdens, and are unsuitable for scenarios when sampling from the barycenter measure is needed.

Given the aforementioned challenges, our work contributes to the literature of "free-support" approaches which do not prescribe any discrete support when approximating the W_2 -barycenter, and our algorithm works for general continuous non-parametric probability measures. From a high-level perspective, we propose an implementable stochastic counterpart to the prominent theoretic fixed-point framework provided by Álvarez-Esteban, Del Barrio, Cuesta-Albertos, and Matrán [3]. Specifically, Álvarez-Esteban et al. [3] demonstrated that the W_2 -barycenter of absolutely continuous probability measures $\nu_1, \ldots, \nu_K \in \mathcal{P}_{2,ac}(\mathbb{R}^d)$ can be computed via a fixed-point of the operator $G : \mathcal{P}_{2,ac}(\mathbb{R}^d) \to \mathcal{P}_{2,ac}(\mathbb{R}^d)$ defined through the pushforward operation:¹

$$G(\mu) := \left[\sum_{k=1}^{K} w_k T^{\mu}_{\nu_k}\right] \sharp \mu \qquad \forall \mu \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d), \tag{1.3}$$

where $T_{\nu_k}^{\mu}$ corresponds to Monge's optimal transport (OT) map from μ to ν_k , i.e.,

$$T^{\mu}_{\nu_{k}} \in \operatorname*{arg\,min}_{T} \left\{ \int_{\mathbb{R}^{d}} \left\| \boldsymbol{x} - T(\boldsymbol{x}) \right\|^{2} \mu(\mathrm{d}\boldsymbol{x}) : T : \mathbb{R}^{d} \to \mathbb{R}^{d} \text{ is Borel measurable and } T \sharp \mu = \nu_{k} \right\}$$

In particular, they showed that the G-operator in (1.3) is continuous with respect to the W_2 -metric [3, Theorem 3.1], and that the following theorem holds.

Theorem 1.2 (Properties of the *G*-operator [3, Corollary 3.5 & Theorem 3.6]). Let $\nu_1, \ldots, \nu_K \in \mathcal{P}_{2,ac}(\mathbb{R}^d)$. The *G*-operator defined in (1.3) satisfies the following properties.

- (i) The unique W_2 -barycenter $\bar{\mu} \in \mathcal{P}_{2,ac}(\mathbb{R}^d)$ (see Theorem 2.3) of ν_1, \ldots, ν_K with weights w_1, \ldots, w_K is a fixed-point of G, i.e., $\bar{\mu} = G(\bar{\mu})$.
- (ii) For any $\mu_0 \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$, the sequence $(\mu_t)_{t \in \mathbb{N}_0}$ generated by the iteration

$$\mu_{t+1} := G(\mu_t) \qquad \forall t \in \mathbb{N}_0 \tag{1.4}$$

is tight. Moreover, every accumulation point of the sequence $(\mu_t)_{t\in\mathbb{N}_0}$ with respect to the \mathcal{W}_2 -metric is a fixed-point of G.

Theorem 1.2 then leads to a simple iterative scheme for W_2 -barycenter where one begins with an arbitrary $\mu_0 \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$ and iterates (1.4) to generate $(\mu_t)_{t\in\mathbb{N}_0}$, and guarantees that $(\mu_t)_{t\in\mathbb{N}_0}$ converges in W_2 to the W_2 -barycenter of ν_1, \ldots, ν_K with weights w_1, \ldots, w_K when G has a unique fixed-point. More recently, it has been shown by Tanguy, Delon, and Gozlan [63] that this fixed-point method can be generalized to compute barycenters under diverse transportation costs and generic measures.

However, when ν_1, \ldots, ν_K are general non-parametric probability measures, the operation (1.3) is a theoretical but impractical "oracle" due to the difficulty in computing the OT map $T^{\mu}_{\nu_k}$ exactly. Therefore, numerical implementations of this scheme are either limited to particular parametric measures from the same elliptical family (see, e.g., [3, Section 4]), or carried out via neural network approximations [38] at the price of analytical difficulties. This bottleneck motivated our development of a rigorous and provably convergent estimator-based

¹For two closed subsets \mathcal{X}, \mathcal{Y} of Euclidean spaces and a Borel measurable function $T : \mathcal{X} \to \mathcal{Y}$, the pushforward of a probability measure $\mu \in \mathcal{P}(\mathcal{X})$ by T is denoted by $T \sharp \mu \in \mathcal{P}(\mathcal{Y})$, which is defined via $T \sharp \mu(B) \equiv \mu \circ T^{-1}(B)$ for every Borel set $B \subseteq \mathcal{Y}$.

3

Conceptual Algorithm 1: Stochastic fixed-point iterative scheme.²

Input: $K \in \mathbb{N}$ input probability measures $\nu_1, \ldots, \nu_K \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$, weights $w_1 > 0, \ldots, w_K > 0$ with $\sum_{k=1}^{K} w_k = 1$, initial probability measure $\mu_0 \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$. **Output:** $(\widehat{\mu}_t)_{t \in \mathbb{N}_0}$. 1 Initialize $\widehat{\mu}_0 \leftarrow \mu_0$. **2** for $t = 0, 1, 2, \dots$ do for $k = 1, \ldots, K$ do 3 Randomly generate $\widehat{N}_{t,k} \in \mathbb{N}$ independent samples $\{X_{t+1,k,i}\}_{i=1:\widehat{N}_{t,k}}$ from $\widehat{\mu}_t$. 4 Randomly generate $\widehat{N}_{t,k} \in \mathbb{N}$ independent samples $\{Y_{t+1,k,i}\}_{i=1:\widehat{N}_{t,k}}$ from ν_k . 5 Approximate $T_{\nu_k}^{\hat{\mu}_t}$ with an estimator $\hat{T}_{t+1,k} \approx T_{\nu_k}^{\hat{\mu}_t}$ using the samples $\{X_{t+1,k,i}\}_{i=1:\hat{N}_{t,k}}$ and 6 $\{\boldsymbol{Y}_{t+1,k,i}\}_{i=1:\widehat{N}_{t,k}}.$ Choose $\widehat{\mu}_{t+1} \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$ such that $\widehat{\mu}_{t+1} \approx \left[\sum_{k=1}^K w_k \widehat{T}_{t+1,k}\right] \sharp \widehat{\mu}_t$. 7 8 return $(\widehat{\mu}_t)_{t \in \mathbb{N}_0}$.

stochastic extension of this deterministic fixed-point iterative scheme, which is beyond the capabilities of existing "free-support" algorithms.

The idea of our stochastic fixed-point iterative scheme is sketched in Conceptual Algorithm 1 whose concrete implementation details are deferred to Section 3. Intuitively, our algorithm generates a sequence $(\hat{\mu}_t)_{t \in \mathbb{N}_0}$ by approximating each true OT map $T_{\nu_k}^{\hat{\mu}_t}$ with an OT map estimator $\hat{T}_{t+1,k}$ (Line 6), and approximating the *G*-operator defined in (1.3) when updating from $\hat{\mu}_t$ to $\hat{\mu}_{t+1}$ (Line 7). In particular, letting $\hat{T}_{t+1,k} = T_{\nu_k}^{\hat{\mu}_t}$ in Line 6 and letting $\hat{\mu}_{t+1} = \left[\sum_{k=1}^{K} w_k \hat{T}_{t+1,k}\right] \sharp \hat{\mu}_t$ in Line 7 will recover the deterministic fixed-point iterative scheme. Our objective is to develop a concrete setting as well as a computationally tractable implementation of Conceptual Algorithm 1 such that the resultant stochastic sequence of probability measures $(\hat{\mu}_t)_{t\in\mathbb{N}_0}$ will converge to the \mathcal{W}_2 -barycenter of ν_1, \ldots, ν_K with weights w_1, \ldots, w_K in an almost sure sense. Specifically, our contributions in this paper can be summarized as follows:

- (i) We provide a computationally tractable stochastic fixed-point algorithm (i.e., Algorithm 2) for approximately computing the W₂-barycenter of general input measures ν₁,..., ν_K. Our iterative algorithm is characterized by a tailored truncating operation and an "admissible" class of OT map estimators (see Assumption 3.4). In particular, we neither restrict the support of the approximate W₂-barycenter to be a finite collection of points nor restrict ν₁,..., ν_K to be discrete or to specific parametric families of measures, with only mild regularity conditions required instead (see Assumption 3.1 and Setting 3.6).
- (ii) We perform a rigorous convergence analysis of our algorithm to show that it converges to the true W_2 barycenter of ν_1, \ldots, ν_K in an almost sure sense (see Setting 3.13 and Theorem 3.14). Concretely, we adapt the computationally efficient entropic OT map estimator designed by Pooladian and Niles-Weed [52] in our algorithm up to tailored modifications to guarantee convergence (see Corollary 4.3). To the best of our knowledge, our algorithm is the first computationally tractable extension of the fixed-point iterative scheme by Álvarez-Esteban et al. [3] with convergence guarantee.
- (iii) We propose a simple and efficient method for generating synthetic instances of the W_2 -barycenter problems where the input measures are continuous and non-parametric and the true W_2 -barycenter is known (see Proposition 5.2). These problem instances can be used to evaluate and compare the effectiveness of W_2 -barycenter algorithms. We demonstrate via numerical experiments that our algorithm is empirically accurate, efficient, and stable compared with state-of-the-art algorithms, and allows implementations in a distributed and parallel computing environment, which is attractive to large-scale application scenarios.

Our paper is organized as follows. Section 1.1 provides a literature review on free-support algorithms for approximating the W_2 -barycenter and on OT map estimators, and Section 1.2 introduces the notations used in this paper. Section 2 mentions the key preliminary results on which our arguments are mainly based. Section 3

²We leave this conceptual algorithm abstract here for ease of illustration. Concrete choices of the sample size (Line 4 and 5), the OT map estimator (Line 6), and the approximate pushforward measure (Line 7) will be specified in Section 3.

presents our stochastic fixed-point algorithm for approximating the W_2 -barycenter, where we also perform detailed analysis of its convergence. In Section 4, we concretely develop a modified entropic OT map estimator that can guarantee the convergence of our stochastic fixed-point algorithm. Section 5 presents our novel algorithm for generating synthetic instances of the W_2 -barycenter problems. Finally, in Section 6, we compare our proposed algorithm with other state-of-the-art methods for W_2 -barycenter approximation based on generative models.

1.1. **Literature review.** In this subsection, we review in detail two streams of research that are closely related to our study, namely free-support methods for approximating Wasserstein barycenters and statistical estimation approaches of the optimal transport map.

Free-support methods for approximating the Wasserstein barycenter. Free-support methods do not anticipate potential supports of the underlying Wasserstein barycenter a priori. As variants to fixed-support schemes, there have been practices on alternating between optimizing the histogram weights and optimizing the support atoms via stochastic optimization [16] or incrementally updating supports via the Frank–Wolfe algorithm [41]. However, these candidates are still subject to the inherent limitations in scalability arising from discrete supports. The last few years have witnessed an extensive and rapid development of algorithms which impose no restrictions on the support of the underlying Wasserstein barycenter, along with the thriving of generative models. These algorithms can be methodologically distinguished into three classes. The first class directly solve the variational problem (1.1) over measures, which are most times bottlenecked by the computational challenge of evaluating the W_2 -distance. For instance, Cohen, Arbel, and Deisenroth [17] detoured to solving the Sinkhorn barycenter as a proxy of the W_2 -barycenter, and Fan, Taghvaei, and Chen [26] parametrized the W_2 -barycenter using generative neural networks and solved an Input Convex Neural Network (ICNN) [5] based min-max-min problem. The second class of algorithms characterize the Brenier potentials (see Theorem 2.4) by solving the dual of (1.1) using reproducing kernel Hilbert space or neural networks [29, 40]. To recover the barycenter from individual potentials, Li et al. [40] considered the barycentric projection [4, Definition 5.4.2] approach while Korotin, Li, Solomon, and Burnaev [37] considered pushforwards by gradients of the potentials. The last class of free-support algorithms initiate from the fixed-point iterative framework derived by Álvarez-Esteban et al. [3], which has been found generalizable to more generic transportation costs [63]. The idea has been numerically implemented by Korotin et al. [38] via generative neural networks, and by von Lindheim [70] via barycentric projection in settings with discrete measures.

Estimation methods of the optimal transport map. Our stochastic fixed-point algorithm consists of estimators of the optimal transport (OT) map between measures. Besides the optimal transportation cost, the OT map itself has been of primary interest in diverse applications including transfer learning, computational biology, nonparametric hypothesis testing, and so on; see, e.g., [14, Chapter 3] and the references therein for a review. However, computing the true OT map is exceptionally hard provided the difficulty of evaluating the Wasserstein distance; see, e.g., [62, 65]. Recently, diverse types of OT map estimators with rigorous statistical guarantees have been proposed. Under pre-specified regularity assumptions, Hütter and Rigollet [34] and Gunsilius [32] both established a $\mathcal{L}^2(\mu)$ -convergence rate: the former proposed a near-optimal OT map estimator via truncated wavelet approximation and the latter obtained from kernel density estimations an upper bound on the $\mathcal{L}^2(\mu)$ risk. Subsequently, Pooladian and Niles-Weed [52] and Deb, Ghosal, and Sen [22] derived OT map estimators via barycentric projection techniques in regularized and non-regularized settings. Manole, Balakrishnan, Niles-Weed, and Wasserman [42] sharpened the upper bound risk in [32] and introduced in addition the so-called "plug-in" estimators, which are built upon the optimal transport plan between the empirical counterparts of measures. Moreover, Vacher, Muzellec, Rudi, Bach, and Vialard [67] and Muzellec, Vacher, Bach, Vialard, and Rudi [46] proposed estimators with comparable convergence rates by employing kernel sum-of-squares [8, Chapter 3] as building blocks. Another remarkable stream of works developed a class of plug-in estimators via smooth and strongly convex regression and interpolation, in light of the underlying functional form of Brenier potential. For instance, Paty, d'Aspremont, and Cuturi [49] formulated a quadratically constrained quadratic program (QCQP) for approximating the Brenier potential leveraging the convex interpolability framework developed by Taylor [66]; Curmei and Hall [19] parametrized the underlying Brenier potentials as polynomials and solved a shape-constrained polynomial regression problem to approximately recover the OT map; González-Sanz, De Lara, Béthune, and Loubes [31] deployed state-of-the-art Lipschitz-constrained generative adversarial networks (GAN) for OT map estimation in regression.

1.2. Notations. In the following, we introduce the terminologies and notations that are used throughout this paper. All vectors are assumed to be column vectors and are denoted by boldface symbols. In particular, for $k \in \mathbb{N}$, $\mathbf{0}_k$ denotes the vector in \mathbb{R}^k with all entries equal to zero. We also use **0** when the dimension can be inferred from the context. We denote by $\langle \cdot, \cdot \rangle$ the Euclidean dot product, i.e., $\langle x, y \rangle := x^{\mathsf{T}}y$ and we denote by $\|\cdot\|$ the Euclidean norm, i.e., $\|x\| := (\langle x, x \rangle)^{\frac{1}{2}}$. Open and closed balls centered at x with radius r are denoted by $B(\boldsymbol{x},r)$ and $\bar{B}(\boldsymbol{x},r)$, respectively. For any set $\mathcal{X} \subseteq \mathbb{R}^k$, we let $int(\mathcal{X})$, $relint(\mathcal{X})$, $rel(\mathcal{X})$, $bd(\mathcal{X})$, and $\operatorname{relbd}(\mathcal{X})$ denote its interior, relative interior, closure, boundary, and relative boundary, respectively. Moreover, for $k \in \mathbb{N}$, let O_k denote the k-by-k zero matrix, and let I_k denote the k-by-k identity matrix. Let \mathbb{S}^k , \mathbb{S}^k_+ , and \mathbb{S}_{++}^k denote the set of k-by-k matrices that are symmetric, symmetric positive semi-definite, and symmetric positive definite, respectively. For $\mathbf{A}, \mathbf{B} \in \mathbb{S}^k$, let $\mathbf{A} \succeq \mathbf{B}$ be equivalent to $\mathbf{A} - \mathbf{B} \in \mathbb{S}^k_+$. Furthermore, the smallest and the largest eigenvalues of any $\mathbf{A} \in \mathbb{S}^k$ are denoted by $e_{\min}(\mathbf{A})$ and $e_{\max}(\mathbf{A})$ respectively.

For a closed subset \mathcal{X} of a Euclidean space, let $\mathcal{B}(\mathcal{X})$ denote the Borel subsets of \mathcal{X} , and let $\mathcal{P}(\mathcal{X})$ denote the set of Borel probability measures on \mathcal{X} , while $\mathcal{P}_2(\mathcal{X}) \subseteq \mathcal{P}(\mathcal{X})$ consists of the ones with finite second moments. As mentioned in Section 1, the associated set $\mathcal{P}_{2,ac}(\mathcal{X})$ contains probability measures in $\mathcal{P}_2(\mathcal{X})$ which are absolutely continuous with respect to the Lebesgue measure. For any $\mu \in \mathcal{P}(\mathcal{X})$ and any $\mathcal{Y} \in \mathcal{B}(\mathcal{X})$ with $\mu(\mathcal{Y}) > 0$, let $\mu|_{\mathcal{Y}}$ denote the probability measure formed by truncating μ to \mathcal{Y} , i.e., $\mu|_{\mathcal{Y}}(A) := \frac{\mu(\mathcal{Y} \cap A)}{\mu(\mathcal{Y})}$ for all $A \in \mathcal{B}(\mathcal{X})$. As mentioned, for closed subsets \mathcal{X}, \mathcal{Y} of Euclidean spaces, the pushforward of a probability measure $\mu \in \mathcal{P}(\mathcal{X})$ by a Borel measurable function $T : \mathcal{X} \to \mathcal{Y}$ is denoted by $T \sharp \mu \in \mathcal{P}(\mathcal{Y})$.

Let us also introduce the notations for the following function classes. For an open bounded set $\mathcal{X} \subset \mathbb{R}^d$ and for $q \in \mathbb{N}_0$, $\alpha \in (0, 1]$, let $\mathcal{C}^q(cl(\mathcal{X}))$ denote the set of \mathbb{R} -valued continuous functions on $cl(\mathcal{X})$ that are q-times continuously differentiable on \mathcal{X} , and let $\mathcal{C}^{q,\alpha}(\operatorname{cl}(\mathcal{X}))$ denote the set of \mathbb{R} -valued continuous functions on $\operatorname{cl}(\mathcal{X})$ that are q-times continuously differentiable on \mathcal{X} whose q-th order partial derivatives are α -Hölder continuous. In particular, $C^{q,\alpha}(cl(\mathcal{X}))$ is a Banach space with respect to the following norm (see, e.g., [25, Theorem 5.1.1]):

$$\|\varphi\|_{\mathcal{C}^{q,\alpha}(\mathrm{cl}(\mathcal{X}))} := \max_{|\beta| \le q} \sup_{\boldsymbol{x} \in \mathcal{X}} \left\{ \left| \partial^{\beta} \varphi(\boldsymbol{x}) \right| \right\} + \max_{|\beta| = q} \sup_{\boldsymbol{x}, \boldsymbol{y} \in \mathcal{X}} \left\{ \frac{\left| \partial^{\beta} \varphi(\boldsymbol{x}) - \partial^{\beta} \varphi(\boldsymbol{y}) \right|}{\|\boldsymbol{x} - \boldsymbol{y}\|^{\alpha}} \right\} \qquad \forall \varphi \in \mathcal{C}^{q,\alpha}(\mathrm{cl}(\mathcal{X})),$$

where $\boldsymbol{\beta} = (\beta_1, \dots, \beta_d) \in \mathbb{N}_0^d$ is a multi-index, $|\boldsymbol{\beta}| := \beta_1 + \dots + \beta_d$, and $\partial^{\boldsymbol{\beta}} \varphi := \frac{\partial^{|\boldsymbol{\beta}|} \varphi}{\partial x_1^{\beta_1} \cdots \partial x_d^{\beta_d}}$ denotes the partial derivative of φ with respect to the multi-index β . We call $\mathcal{C}^{q,\alpha}(\operatorname{cl}(\mathcal{X}))$ the set of (q, α) -Hölder functions on $cl(\mathcal{X})$. Moreover, let $\mathcal{C}^{loc,q,\alpha}(\mathbb{R}^d)$ denote the set of \mathbb{R} -valued functions on \mathbb{R}^d that are (q,α) -Hölder when restricted to the closure of any bounded open set. We call $\mathcal{C}^{\text{loc},q,\alpha}(\mathbb{R}^d)$ the set of locally (q,α) -Hölder functions on \mathbb{R}^d . In addition, let $\mathcal{C}^{\infty}(\mathbb{R}^d)$ denote the set of infinitely differentiable \mathbb{R} -valued functions on \mathbb{R}^d . Lastly, we denote by $\mathcal{C}_{\text{lin}}(\mathbb{R}^d, \mathbb{R}^d)$ the set of continuous functions from \mathbb{R}^d to \mathbb{R}^d that have at most linear growth, i.e., $T \in$ $\mathcal{C}_{\text{lin}}(\mathbb{R}^{d},\mathbb{R}^{d}) \text{ if and only if } T:\mathbb{R}^{d} \to \mathbb{R}^{d} \text{ is continuous and } \sup_{\boldsymbol{x}\in\mathbb{R}^{d}} \left\{\frac{\|T(\boldsymbol{x})\|}{1+\|\boldsymbol{x}\|}\right\} < \infty. \text{ Note that } \mathcal{C}_{\text{lin}}(\mathbb{R}^{d},\mathbb{R}^{d}) \text{ is a Banach space with respect to the norm } \|T\|_{\mathcal{C}_{\text{lin}}(\mathbb{R}^{d},\mathbb{R}^{d})} := \sup_{\boldsymbol{x}\in\mathbb{R}^{d}} \left\{\frac{\|T(\boldsymbol{x})\|}{1+\|\boldsymbol{x}\|}\right\} \forall T \in \mathcal{C}_{\text{lin}}(\mathbb{R}^{d},\mathbb{R}^{d}).$

Furthermore, we denote by $\mathfrak{C}_{\lambda,\overline{\lambda}}(\mathbb{R}^d)$ the collection of proper, lower semi-continuous (l.s.c.), and convex functions on \mathbb{R}^d which are $\overline{\lambda}$ -smooth and $\underline{\lambda}$ -strongly convex; see Definition 2.1. In particular, $\mathfrak{C}_{0,\infty}(\mathbb{R}^d)$ contains all proper l.s.c. convex functions on \mathbb{R}^d . The subdifferential of any $\varphi \in \mathfrak{C}_{0,\infty}(\mathbb{R}^d)$ at $\boldsymbol{x} \in \mathbb{R}^d$ is denoted by $\partial \varphi(\boldsymbol{x})$. In addition, we denote $\mathfrak{C}^{\infty}_{\underline{\lambda},\overline{\lambda}}(\mathbb{R}^d) \cap \mathfrak{C}_{\underline{\lambda},\overline{\lambda}}(\mathbb{R}^d), \mathfrak{C}^q_{\underline{\lambda},\overline{\lambda}}(\mathbb{R}^d) := \mathcal{C}^q(\mathbb{R}^d) \cap \mathfrak{C}_{\underline{\lambda},\overline{\lambda}}(\mathbb{R}^d)$, $\mathfrak{C}^{\mathrm{loc},q,\alpha}_{\lambda,\overline{\lambda}}(\mathbb{R}^d) := \mathcal{C}^{\mathrm{loc},q,\alpha}(\mathbb{R}^d) \cap \mathfrak{C}_{\underline{\lambda},\overline{\lambda}}(\mathbb{R}^d) \text{ for } q \in \mathbb{N}_0, \alpha \in (0,1].$

2. PRELIMINARY RESULTS

In this section, we provide an overview of the preliminary results that are frequently used in our discussions. Readers who are familiar with the optimal transport theory can skip this part and proceed to Section 3.

The notions of smooth and strongly convex functions are formally defined as follows.

Definition 2.1 (Smooth and strongly convex functions). For $0 \le \underline{\lambda} \le \overline{\lambda} \le \infty$, a proper, lower semi-continuous (l.s.c.) and convex function $\varphi : \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ is called $\overline{\lambda}$ -smooth (i.e., $\varphi \in \mathfrak{C}_{0,\overline{\lambda}}(\mathbb{R}^d)$) if

$$arphi(oldsymbol{y}) \leq arphi(oldsymbol{x}) + \langle oldsymbol{g}, oldsymbol{y} - oldsymbol{x}
angle + rac{\overline{\lambda}}{2} \|oldsymbol{x} - oldsymbol{y}\|^2 \qquad orall oldsymbol{x}, oldsymbol{y} \in \mathbb{R}^d, \ orall oldsymbol{g} \in \partial arphi(oldsymbol{x}),$$

and is called $\underline{\lambda}$ -strongly convex (i.e., $\varphi \in \mathfrak{C}_{\lambda,\infty}(\mathbb{R}^d)$) if

$$arphi(oldsymbol{y}) \geq arphi(oldsymbol{x}) + \langle oldsymbol{g}, oldsymbol{y} - oldsymbol{x}
angle + rac{\lambda}{2} \|oldsymbol{x} - oldsymbol{y}\|^2 \qquad orall oldsymbol{x}, oldsymbol{y} \in \mathbb{R}^d, \ orall oldsymbol{g} \in \partial arphi(oldsymbol{x})$$

It follows from classical results (see, e.g., [47, Lemma 1.2.3 & Theorem 2.1.5]) that for $\overline{\lambda} < \infty$, every $\varphi \in \mathfrak{C}_{0,\overline{\lambda}}(\mathbb{R}^d)$ is continuously differentiable on \mathbb{R}^d and $\nabla \varphi$ is $\overline{\lambda}$ -Lipschitz continuous.

Many of our discussions in this paper invoke results from the optimal transport theory and properties around the Wasserstein distance between probability measures; see, e.g., the books of Villani [68, 69] and Santambrogio [58]. We start by introducing the notion of couplings.

Definition 2.2 (Coupling). Given $m \in \mathbb{N}$ probability measures $\nu_1 \in \mathcal{P}(\mathcal{X}_1), \ldots, \nu_m \in \mathcal{P}(\mathcal{X}_m)$ on closed subsets $\mathcal{X}_1, \ldots, \mathcal{X}_m$ of \mathbb{R}^d , the set of couplings of ν_1, \ldots, ν_m is denoted by $\Pi(\nu_1, \ldots, \nu_m)$, which is defined as

$$\Pi(\nu_1,\ldots,\nu_m) := \{ \pi \in \mathcal{P}(\mathcal{X}_1 \times \cdots \times \mathcal{X}_m) : \text{ the marginal of } \pi \text{ on } \mathcal{X}_i \text{ is } \nu_i \text{ for } i = 1,\ldots,m \}.$$

The minimization problem embedded in the formulation (1.1) is known as Kantorovich's optimal transport problem [35] with respect to the squared-distance cost, and the infimum is well known to be attained by an optimal coupling; see, e.g., [69, Theorem 4.1]. In the rest of this paper, the optimality of a coupling is always considered with respect to the squared-distance cost. The existence of a W_2 -barycenter is shown by Agueh and Carlier [1, Proposition 2.3], and there may exist more than one W_2 -barycenters of $\nu_1, \nu_2, \ldots, \nu_K$ in general. A sufficient condition to guarantee the uniqueness of W_2 -barycenter is given as follows.

Theorem 2.3 ([1, Proposition 3.5 & Theorem 5.1]). Among $\nu_1, \ldots, \nu_K \in \mathcal{P}_2(\mathbb{R}^d)$, if there exists at least one index $k \in \{1, \ldots, K\}$ such that $\nu_k \in \mathcal{P}_{2,ac}(\mathbb{R}^d)$, then the \mathcal{W}_2 -barycenter $\bar{\mu}$ in Definition 1.1 is unique. Moreover, if there exists at least one index $k \in \{1, \ldots, K\}$ such that ν_k has bounded density, then the unique $\bar{\mu} \in \mathcal{P}_{2,ac}(\mathbb{R}^d)$.

Next, let us present Brenier's theorem which characterizes optimal couplings with gradient of convex functions when the source measure μ is absolutely continuous with respect to the Lebesgue measure; see, e.g., [68, Theorem 2.12].

Theorem 2.4 (Brenier's theorem). Let $\mu \in \mathcal{P}_{2,ac}(\mathbb{R}^d)$, $\nu \in \mathcal{P}_2(\mathbb{R}^d)$. Then, there is a unique optimal coupling $\pi^* \in \Pi(\mu, \nu)$ that minimizes (1.1). Moreover, $\pi \in \Pi(\mu, \nu)$ minimizes (1.1) if and only if there exists a proper, *l.s.c.*, and convex function $\varphi_{\nu}^{\mu} : \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ such that $\pi = [I_d, T_{\nu}^{\mu}] \sharp \mu$ where $I_d : \mathbb{R}^d \to \mathbb{R}^d$ denotes the identity map on \mathbb{R}^d and $T_{\nu}^{\mu} = \nabla \varphi_{\nu}^{\mu}$ is the μ -a.e. everywhere unique gradient of φ_{ν}^{μ} . In this case, it holds that

$$\mathcal{W}_2(\mu,\nu)^2 = \int_{\mathbb{R}^d} \|\boldsymbol{x}\|^2 - 2\varphi_{\nu}^{\mu}(\boldsymbol{x})\,\mu(\mathrm{d}\boldsymbol{x}) + \int_{\mathbb{R}^d} \|\boldsymbol{y}\|^2 - 2\sup_{\boldsymbol{x}\in\mathbb{R}^d} \left\{ \langle \boldsymbol{y}, \boldsymbol{x} \rangle - \varphi_{\nu}^{\mu}(\boldsymbol{x}) \right\} \nu(\mathrm{d}\boldsymbol{y}),$$

and T^{μ}_{ν} is the μ -almost everywhere unique optimal solution of Monge's optimal transport problem:

$$\inf\left\{\int_{\mathbb{R}^d} \left\|\boldsymbol{x} - T(\boldsymbol{x})\right\|^2 \mu(\mathrm{d}\boldsymbol{x}) : T : \mathbb{R}^d \to \mathbb{R}^d \text{ is Borel measurable and } T \sharp \mu = \nu\right\}$$

We refer to $\varphi_{\nu}^{\mu} : \mathbb{R}^{d} \to \mathbb{R} \cup \{\infty\}$ and $T_{\nu}^{\mu} : \mathbb{R}^{d} \to \mathbb{R}^{d}$ in Theorem 2.4 as the Brenier potential from μ to ν and the optimal transport (OT) map from μ to ν , respectively. In general, the μ -almost everywhere uniqueness of T_{ν}^{μ} does not necessarily imply the μ -almost everywhere uniqueness of φ_{ν}^{μ} even up to an additive constant. However, φ_{ν}^{μ} becomes μ -almost everywhere uniquely determined up to an additive constant if $\operatorname{supp}(\mu)$ is the closure of a connected open set on which μ has positive density; see, e.g., [69, Remark 10.30].

Furthermore, the convergence of our proposed algorithm requires regularity properties of the OT map T_{ν}^{μ} . Regarding this matter, a series of studies by Caffarelli [9, 10, 11, 12] developed the foundations of the regularity theory of OT maps under suitable geometric assumptions on the supports and densities of the measures. Here, we partially report these results as phrased in [69, Theorem 12.50].

Theorem 2.5 (Caffarelli's global regularity theory). Let \mathcal{X}_{μ} and \mathcal{X}_{ν} be two connected bounded open sets in \mathbb{R}^d that both have \mathcal{C}^2 -boundaries and are both uniformly convex.³ Let $\mu \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$ be concentrated on \mathcal{X}_{μ} , and let $\nu \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$ be concentrated on \mathcal{X}_{ν} , i.e., $\mu(\mathbb{R}^d \setminus \mathcal{X}_{\mu}) = \nu(\mathbb{R}^d \setminus \mathcal{X}_{\nu}) = 0$. Suppose that for $q \in \mathbb{N}_0$,

³A set $\mathcal{X} \subset \mathbb{R}^d$ is said to have \mathcal{C}^p boundary with $p \in [0, \infty)$ if $\operatorname{bd}(\mathcal{X})$ is locally the graph of a \mathcal{C}^p function, and is said to be uniformly convex if for every $\epsilon > 0$, there exists $\delta > 0$ such that, for any $\boldsymbol{x}, \boldsymbol{y} \in \mathcal{X}$ with $\|\boldsymbol{x} - \boldsymbol{y}\| < \epsilon$, the distance from the mid-point $(\boldsymbol{x} + \boldsymbol{y})/2$ to $\operatorname{bd}(\mathcal{X})$ is at least δ .

 $\alpha \in (0,1], f_{\mu} \in C^{q,\alpha}(\operatorname{cl}(\mathcal{X}_{\mu})) \text{ and } f_{\nu} \in C^{q,\alpha}(\operatorname{cl}(\mathcal{X}_{\nu})) \text{ are the density functions of } \mu \text{ and } \nu \text{ with respect to the Lebesgue measure, respectively. Moreover, suppose that there exists } \gamma > 1 \text{ such that } \gamma^{-1} \leq f_{\mu}(\boldsymbol{x}) \leq \gamma \text{ for all } \boldsymbol{x} \in \operatorname{cl}(\mathcal{X}_{\mu}) \text{ and that } \gamma^{-1} \leq f_{\nu}(\boldsymbol{x}) \leq \gamma \text{ for all } \boldsymbol{x} \in \operatorname{cl}(\mathcal{X}_{\nu}). \text{ Then, the Brenier potential } \varphi_{\nu}^{\mu} \text{ satisfies } \varphi_{\nu}^{\mu} \in C^{q+2,\alpha}(\operatorname{cl}(\mathcal{X}_{\mu})).$

3. Stochastic fixed-point algorithm for \mathcal{W}_2 -barycenter

In this section, we will present our computationally tractable stochastic fixed-point algorithm for W_2 barycenter and show its convergence. Section 3.1 introduces the specifications of the approximation steps in Line 6 and Line 7 of Conceptual Algorithm 1 as well as additional assumptions. In Section 3.2, we develop sufficient conditions for the convergence of our stochastic fixed-point algorithm.

3.1. **Settings.** Conceptual Algorithm 1 has illustrated the conceptual procedure of our stochastic fixed-point iterative scheme. Before presenting its computationally tractable implementation as a concrete algorithm, let us introduce some additional notions in Definition 3.1, Assumption 3.3, and Assumption 3.4.

Definition 3.1 (Admissible support sets and admissible probability measures). For $d \in \mathbb{N}$, let $S(\mathbb{R}^d)$ denote the collection of subsets of \mathbb{R}^d defined as follows:

 $\mathcal{S}(\mathbb{R}^d) := \{ cl(\mathcal{Y}) : \mathcal{Y} \subset \mathbb{R}^d \text{ is non-empty, open, bounded, uniformly convex, and has a } \mathcal{C}^2\text{-boundary} \}.$

We will refer to $S(\mathbb{R}^d)$ as the admissible support sets. For $q \in \mathbb{N}_0$, let $\mathcal{M}^q(\mathbb{R}^d)$ denote the collection of probability measures on \mathbb{R}^d defined as follows:

$$\mathcal{M}^{q}(\mathbb{R}^{d}) := \left\{ \mu \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^{d}) : \sup_{\gamma^{-1} \leq f_{\mu}(\boldsymbol{x}) \leq \gamma \; \forall \boldsymbol{x} \in \mathrm{supp}(\mu), \; f_{\mu} \text{ is the density function of } \mu \right\}$$

We will refer to $\mathcal{M}^q(\mathbb{R}^d)$ as the set of q-admissible compactly supported probability measures. Moreover, let $\mathcal{M}^q_{full}(\mathbb{R}^d)$ denote the collection of probability measures on \mathbb{R}^d defined as follows:

$$\mathcal{M}^{q}_{\text{full}}(\mathbb{R}^{d}) := \left\{ \rho \in \mathcal{P}_{2,\text{ac}}(\mathbb{R}^{d}) : \sup_{f_{\rho}(\boldsymbol{x}) > 0} \forall \boldsymbol{x} \in \mathbb{R}^{d}, \ f_{\rho} \text{ is the density function of } \rho \right\}.$$

We will refer to $\mathcal{M}^q_{\text{full}}(\mathbb{R}^d)$ as the set of q-admissible fully supported probability measures.

The conditions in the definitions of the admissible support sets and the q-admissible compactly supported probability measures are motivated by the conditions in Caffarelli's global regularity theory (Theorem 2.5). As a result, one can derive the following curvature properties of the Brenier potential φ_{ν}^{μ} from $\mu \in \mathcal{M}^q(\mathbb{R}^d)$ to $\nu \in \mathcal{M}^q(\mathbb{R}^d)$; see, e.g., [42, Lemma 2] and [30, Corollary 3.2]. We highlight that such curvature properties will serve as crucial premises when we control the estimation errors of OT map estimators; see details in Section 4.

Lemma 3.2 (Curvature properties of φ_{ν}^{μ} ; see, e.g., [42, Lemma 2] & [30, Corollary 3.2]). Let $q \in \mathbb{N}_0$, let $\mu, \nu \in \mathcal{M}^q(\mathbb{R}^d)$ be arbitrary, and let $\varphi_{\nu}^{\mu} : \mathbb{R}^d \to \mathbb{R}$ be the Brenier potential from μ to ν (that is unique μ -almost everywhere up to the addition of an arbitrary constant by [69, Remark 10.30]). Then, $\varphi_{\nu}^{\mu} \in \mathcal{C}^{q+2}(\operatorname{supp}(\mu))$ and there exist $0 < \lambda_{\text{LB}} \le \lambda_{\text{UB}} < \infty$ such that $\lambda_{\text{LB}}\mathbf{I}_d \preceq \nabla^2 \varphi_{\nu}^{\mu}(\mathbf{x}) \preceq \lambda_{\text{UB}}\mathbf{I}_d$ for all $\mathbf{x} \in \operatorname{supp}(\mu)$. Moreover, there exists $\tilde{\varphi}_{\nu}^{\mu} \in \mathfrak{C}_{\lambda_{\text{LB}},\lambda_{\text{UB}}}(\mathbb{R}^d)$ that is equal to φ_{ν}^{μ} on $\operatorname{supp}(\mu)$, and therefore one can let $\varphi_{\nu}^{\mu} \in \mathfrak{C}_{\lambda_{\text{LB}},\lambda_{\text{UB}}}(\mathbb{R}^d)$ without loss of generality.

Proof of Lemma 3.2. Let f_{μ} and f_{ν} denote the density functions of μ and ν which satisfy the conditions in Definition 3.1. Thus, $f_{\mu} \in C^{q,\alpha}(\operatorname{supp}(\mu))$ and $f_{\nu} \in C^{q,\alpha'}(\operatorname{supp}(\nu))$ for some $\alpha, \alpha' \in (0, 1]$. This implies that $f_{\mu} \in C^{q,\alpha''}(\operatorname{supp}(\mu))$ and $f_{\nu} \in C^{q,\alpha''}(\operatorname{supp}(\nu))$ for $\alpha'' := \min\{\alpha, \alpha'\}$, and hence Caffarelli's global regularity theory (Theorem 2.5) implies that the Brenier potential φ^{μ}_{ν} satisfies $\varphi^{\mu}_{\nu} \in C^{q+2,\alpha''}(\operatorname{supp}(\mu))$. Thus, the compactness of $\operatorname{supp}(\mu)$ implies that there exists $\lambda_{\mathrm{UB}} < \infty$ such that $\nabla^2 \varphi^{\mu}_{\nu}(x) \preceq \lambda_{\mathrm{UB}} \mathbf{I}_d$ for all $x \in$ $\operatorname{supp}(\mu)$. Moreover, φ^{μ}_{ν} needs to satisfy the following Monge–Ampère type equation as implied by the change of variable formula for pushforward (see, e.g., [4, Lemma 5.5.3]):

$$\det \left(\nabla^2 \varphi^{\mu}_{\nu}(\boldsymbol{x})\right) = \frac{f_{\mu}(\boldsymbol{x})}{f_{\nu} \left(\nabla \varphi^{\mu}_{\nu}(\boldsymbol{x})\right)} \qquad \forall \boldsymbol{x} \in \operatorname{supp}(\mu)$$

Since f_{ν} is bounded from above and f_{μ} is bounded away from 0 on $\operatorname{supp}(\mu)$, it follows that $\det \left(\nabla^2 \varphi_{\nu}^{\mu}(\boldsymbol{x}) \right)$ is bounded away from 0 on $\operatorname{supp}(\mu)$. Combining this and $\nabla^2 \varphi_{\nu}^{\mu}(\boldsymbol{x}) \preceq \lambda_{\operatorname{UB}} \mathbf{I}_d \ \forall \boldsymbol{x} \in \operatorname{supp}(\mu)$ shows that there

exists $\lambda_{\text{LB}} > 0$ such that $\nabla^2 \varphi^{\mu}_{\nu}(\boldsymbol{x}) \succeq \lambda_{\text{LB}} \mathbf{I}_d$ for all $\boldsymbol{x} \in \text{supp}(\mu)$. Consequently, the convexity of $\text{supp}(\mu)$ and the mean value version of Taylor's theorem yield

$$\begin{split} \varphi_{\nu}^{\mu}(\boldsymbol{y}) &\geq \varphi_{\nu}^{\mu}(\boldsymbol{x}) + \langle \nabla \varphi_{\nu}^{\mu}(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle + \frac{\lambda_{\text{LB}}}{2} \|\boldsymbol{x} - \boldsymbol{y}\|^2 & \quad \forall \boldsymbol{x}, \boldsymbol{y} \in \text{supp}(\mu), \\ \varphi_{\nu}^{\mu}(\boldsymbol{y}) &\leq \varphi_{\nu}^{\mu}(\boldsymbol{x}) + \langle \nabla \varphi_{\nu}^{\mu}(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle + \frac{\lambda_{\text{UB}}}{2} \|\boldsymbol{x} - \boldsymbol{y}\|^2 & \quad \forall \boldsymbol{x}, \boldsymbol{y} \in \text{supp}(\mu). \end{split}$$

We can then apply [66, Theorem 2.57 & Remark 2.59] to extend φ^{μ}_{ν} to $\tilde{\varphi}^{\mu}_{\nu} \in \mathfrak{C}_{\lambda_{\mathrm{LB}},\lambda_{\mathrm{UB}}}(\mathbb{R}^d)$ such that $\varphi^{\mu}_{\nu}(\boldsymbol{x}) = \widetilde{\varphi}^{\mu}_{\nu}(\boldsymbol{x}), \nabla \varphi^{\mu}_{\nu}(\boldsymbol{x}) = \nabla \widetilde{\varphi}^{\mu}_{\nu}(\boldsymbol{x})$ for all $\boldsymbol{x} \in \mathrm{supp}(\mu)$. The proof is now complete.

Additionally, our algorithm requires a family of sets on \mathbb{R}^d for truncating probability measures in $\mathcal{M}^q_{\text{full}}(\mathbb{R}^d)$ to probability measures in $\mathcal{M}^q(\mathbb{R}^d)$. The adopted family of sets needs to satisfy the assumption below.

Assumption 3.3 (Family of increasing sets). $(\mathcal{X}_r)_{r\in\mathbb{N}}$ is an infinite collection of subsets of \mathbb{R}^d that satisfies: $\mathcal{X}_r \in \mathcal{S}(\mathbb{R}^d)$, $\mathcal{X}_{r+1} \supseteq \mathcal{X}_r \forall r \in \mathbb{N}$, and $\bigcup_{r\in\mathbb{N}} \mathcal{X}_r = \mathbb{R}^d$.

A concrete example of such a family of increasing sets is $(\bar{B}(\mathbf{0}_d, r))_{r \in \mathbb{N}}$. Similarly, a family of increasing ellipsoids in \mathbb{R}^d also satisfies Assumption 3.3.

Furthermore, with respect to any pair of q-admissible compactly supported probability measures $\mu, \nu \in \mathcal{M}^q(\mathbb{R}^d)$, we consider estimators of the true OT map T^{μ}_{ν} from μ to ν which satisfy the conditions below.

Assumption 3.4 (Admissible OT map estimator). Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space. For $q \in \mathbb{N}_0$, let $\mu, \nu \in \mathcal{M}^q(\mathbb{R}^d)$ and let $\underline{m}, \underline{n} \in \mathbb{N}$ be constants that do not depend on μ or ν . For $m \geq \underline{m}$ and $n \geq \underline{n}$, let $X_1, \ldots, X_m, Y_1, \ldots, Y_n : \Omega \to \mathbb{R}^d$ be independent random variables such that $\operatorname{law}(X_i) = \mu$ for $i = 1, \ldots, m$ and $\operatorname{law}(Y_j) = \nu$ for $j = 1, \ldots, n$. Let Θ be a metric space, where each $\theta \in \Theta$ denotes the parameter(s) that may, for example, represent the extent of smoothing/regularization (see Section 4 for details about the parameter in a concrete OT map estimator). Subsequently, for any $\theta \in \Theta$, let $\widehat{T}_{\nu,n}^{\mu,m}[X_1, \ldots, X_m, Y_1, \ldots, Y_n, \theta] \in C_{\text{lin}}(\mathbb{R}^d, \mathbb{R}^d)$ estimate the OT map $T_{\nu,n}^{\mu}$ from μ to ν based on the samples X_1, \ldots, X_m from μ and the samples Y_1, \ldots, Y_n from ν , where $\widehat{T}_{\nu,n}^{\mu,m}[X_1, \ldots, X_m, Y_1, \ldots, Y_n, \theta]$ has a Borel dependence on $(X_1, \ldots, X_m, Y_1, \ldots, Y_n, \theta) \in \mathbb{R}^{md} \times \mathbb{R}^{nd} \times \Theta$.⁴ For notational simplicity, we often make the dependence of this estimated OT map on the samples implicit and use $\widehat{T}_{\nu,n}^{\mu,m}[\theta](x) \in \mathbb{R}^d$ to denote $\widehat{T}_{\nu,n}^{\mu,m}[X_1, \ldots, X_m, Y_1, \ldots, Y_n, \theta]$ evaluated at $x \in \mathbb{R}^d$.

We assume that $\widehat{T}_{\nu,n}^{\mu,m}[\theta]$ satisfies the following conditions.

- (i) Shape condition: there exist $\alpha(\mu,\nu,m,n, X_1, \ldots, X_m, Y_1, \ldots, Y_n, \theta) \in (0,1]$ and $\underline{\lambda}(\mu,\nu,m,n, X_1, \ldots, X_m, Y_1, \ldots, Y_n, \theta) \in \mathbb{R}_+$, abbreviated to α and $\underline{\lambda}$, both having a Borel dependence on $(\mu,\nu,m,n, X_1, \ldots, X_m, Y_1, \ldots, Y_n, \theta)$, such that $\alpha \in (0,1]$, $\underline{\lambda} > 0$ hold \mathbb{P} -almost surely, and it holds \mathbb{P} -almost surely that $\widehat{T}_{\nu,n}^{\mu,m}[\theta] = \nabla \widehat{\varphi}_{\nu,n}^{\mu,m}[\theta]$ for a function $\widehat{\varphi}_{\nu,n}^{\mu,m}[\theta] \in \mathfrak{C}_{\underline{\lambda},\infty}^{\mathrm{loc},q+2,\alpha}(\mathbb{R}^d)$.
- (ii) Growth condition: there exist $u_0(\nu), u_1(\nu) \in \mathbb{R}_+$ that only depend on ν such that, for all $m \geq \underline{m}, n \geq \underline{n}, \theta \in \Theta$, and all $\boldsymbol{x} \in \mathbb{R}^d$, it holds that $\mathbb{E}\left[\left\|\widehat{T}_{\nu,n}^{\mu,m}[\theta](\boldsymbol{x}) \widehat{T}_{\nu,n}^{\mu,m}[\theta](\boldsymbol{0})\right\|^2\right] \leq u_0(\nu) + u_1(\nu)\|\boldsymbol{x}\|^2.$
- (iii) Consistency condition: for any $\epsilon > 0$, there exist $\overline{n}(\mu, \nu, \epsilon) \in \mathbb{N}$ that has a Borel dependence on (μ, ν, ϵ) and $\tilde{\theta}(\mu, \nu, m, n, \epsilon) \in \Theta$ that has a Borel dependence on $(\mu, \nu, m, n, \epsilon)$ such that $\overline{n}(\mu, \nu, \epsilon) \geq \max\{\underline{m}, \underline{n}\}$ and

$$\mathbb{E}\Big[\big\|\widehat{T}^{\mu,m}_{\nu,n}\big[\widetilde{\theta}(\mu,\nu,m,n,\epsilon)\big] - T^{\mu}_{\nu}\big\|^{2}_{\mathcal{L}^{2}(\mu)}\Big] \leq \epsilon \qquad \forall m \geq \overline{n}(\mu,\nu,\epsilon), \ \forall n \geq \overline{n}(\mu,\nu,\epsilon),$$

where

$$\left\|\widehat{T}_{\nu,n}^{\mu,m}[\theta] - T_{\nu}^{\mu}\right\|_{\mathcal{L}^{2}(\mu)} := \left(\int_{\mathbb{R}^{d}} \left\|\widehat{T}_{\nu,n}^{\mu,m}[\theta](\boldsymbol{x}) - T_{\nu}^{\mu}(\boldsymbol{x})\right\|^{2} \mu(\mathrm{d}\boldsymbol{x})\right)^{\frac{1}{2}} \qquad \forall \theta \in \Theta$$

A concrete example of OT map estimator which satisfies Assumption 3.4 will be introduced later in Section 4. Note that the consistency condition in Assumption 3.4(iii) is possible due to the curvature properties of $T^{\mu}_{\nu} = \nabla \varphi^{\mu}_{\nu}$ in Lemma 3.2. A crucial consequence of the shape condition and the growth condition of the OT map

⁴We say that $\widehat{T}_{\nu,n}^{\mu,m}[\mathbf{X}_1,\ldots,\mathbf{X}_m,\mathbf{Y}_1,\ldots,\mathbf{Y}_n,\theta] \in \mathcal{C}_{\mathrm{lin}}(\mathbb{R}^d,\mathbb{R}^d)$ has a Borel dependence on $(\mathbf{X}_1,\ldots,\mathbf{X}_m,\mathbf{Y}_1,\ldots,\mathbf{Y}_n,\theta)$ if $\widehat{T}_{\nu,n}^{\mu,m}[\cdot]: (\mathbb{R}^d)^m \times (\mathbb{R}^d)^n \times \Theta \to \mathcal{C}_{\mathrm{lin}}(\mathbb{R}^d,\mathbb{R}^d)$ is Borel measurable. Analogous notions of Borel dependency apply to subsequent arguments.

9

estimator is that it allows us to preserve the regularity properties of the pushforward of a probability measure $\rho \in \mathcal{M}^q_{\text{full}}(\mathbb{R}^d)$, which is stated in the following proposition.

Proposition 3.5 (Preservation of regularity properties). Let $q \in \mathbb{N}_0$ be arbitrary. The following statements hold.

- (i) For any $\rho \in \mathcal{M}^q_{\text{full}}(\mathbb{R}^d)$ and $\mathcal{X} \in \mathcal{S}(\mathbb{R}^d)$, it holds that $\rho|_{\mathcal{X}} \in \mathcal{M}^q(\mathbb{R}^d)$. (ii) For $\underline{\lambda} > 0$ and $\varphi \in \mathfrak{C}^2_{\underline{\lambda},\infty}(\mathbb{R}^d)$, it holds that $T := \nabla \varphi : \mathbb{R}^d \to \mathbb{R}^d$ is a homeomorphism. Moreover, if $T \in \mathcal{C}_{\text{lin}}(\mathbb{R}^d, \mathbb{R}^d)$, then $T \not \equiv \rho \in \mathcal{P}_{2,\text{ac}}(\mathbb{R}^d)$ for any $\rho \in \mathcal{P}_{2,\text{ac}}(\mathbb{R}^d)$.
- (iii) Let $\varphi \in \mathfrak{C}^{\mathrm{loc},q+2,\alpha}_{\underline{\lambda},\infty}(\mathbb{R}^d)$ for some $\alpha \in (0,1]$, $\underline{\lambda} > 0$, and suppose that $T := \nabla \varphi \in \mathcal{C}_{\mathrm{lin}}(\mathbb{R}^d,\mathbb{R}^d)$. Then,
- (iii) Let $\varphi \in C_{\underline{\lambda},\infty}$ (iii) for some $u \in (0,1], \underline{n} \neq 0$, and suppose that $T : \Psi \varphi \in C_{\mathrm{III}}(\mathbb{R}^d, \mathbb{R}^d)$, it holds that T is a homeomorphism and $T \sharp \rho \in \mathcal{M}_{\mathrm{full}}^q(\mathbb{R}^d)$ for any $\rho \in \mathcal{M}_{\mathrm{full}}^q(\mathbb{R}^d)$. (iv) Let $\mu, \nu_1, \ldots, \nu_K \in \mathcal{M}^q(\mathbb{R}^d)$, and let $w_1 > 0, \ldots, w_K > 0$ satisfy $\sum_{k=1}^K w_k = 1$. Let $\widehat{T}_{\nu,n}^{\mu,m}[\theta]$ be an OT map estimator that satisfies Assumption 3.4. Subsequently, for $k = 1, \ldots, K$, let $m_k \geq \underline{m}, n_k \geq \underline{n},$ $\theta_k \in \Theta$ be arbitrary ($\underline{m}, \underline{n} \in \mathbb{N}$ are given by Assumption 3.4), and let $\overline{T} := \sum_{k=1}^K w_k \widehat{T}_{\nu_k, n_k}^{\mu,m_k}[\theta_k]$. Then, it holds \mathbb{P} -almost surely that \overline{T} is a homeomorphism and $\overline{T} \sharp \rho \in \mathcal{M}_{\mathrm{full}}^q(\mathbb{R}^d)$ for any $\rho \in \mathcal{M}_{\mathrm{full}}^q(\mathbb{R}^d)$.

Proof of Proposition 3.5. Let us first prove statement (i). Let $\mu := \rho|_{\mathcal{X}}$. Since $\operatorname{supp}(\rho) = \mathbb{R}^d$, it holds that $\operatorname{supp}(\mu) = \operatorname{supp}(\rho|_{\mathcal{X}}) = \mathcal{X} \in \mathcal{S}(\mathbb{R}^d)$. Moreover, Definition 3.1 states that the density function f_ρ of ρ satisfies $f_{\rho}(\boldsymbol{x}) > 0 \ \forall \boldsymbol{x} \in \mathbb{R}^d$ and $f_{\rho} \in \mathcal{C}^{\mathrm{loc},q,\alpha}(\mathbb{R}^d)$ for some $\alpha \in (0,1]$. It thus holds that $f_{\mu} := \frac{f_{\rho}\mathbb{I}_{\mathcal{X}}}{\rho(\mathcal{X})} \in \mathbb{R}^d$ $\mathcal{C}^{q,\alpha}(\mathcal{X})$ is the density function of μ . Subsequently, the compactness of \mathcal{X} implies that $0 < \inf_{\boldsymbol{x} \in \mathcal{X}} \{f_{\mu}(\boldsymbol{x})\} \leq 1$ $\sup_{x \in \mathcal{X}} \{f_{\mu}(x)\} < \infty$ and thus $\mu \in \mathcal{M}^q(\mathbb{R}^d)$. This completes the proof of statement (i).

Next, let us prove statement (ii). It follows from the duality between smooth convex functions and strongly convex functions (see, e.g., [55, Theorem 26.6]) that T is a homeomorphism. Moreover, since φ is twice continuously differentiable on \mathbb{R}^d , it holds by the second-order characterization of strongly convex functions (see, e.g., [47, Theorem 2.1.6]) that $\nabla^2 \varphi(\boldsymbol{x}) \succeq \underline{\lambda} \mathbf{I}_d$ for all $\boldsymbol{x} \in \mathbb{R}^d$. Let us now assume in addition that $T \in \mathcal{C}_{\text{lin}}(\mathbb{R}^d, \mathbb{R}^d)$, fix an arbitrary $\rho \in \mathcal{P}_{2,ac}(\mathbb{R}^d)$, and let f_{ρ} denote the density function of ρ . Subsequently, the change of variable formula for pushforward (see, e.g., [4, Lemma 5.5.3]) yields the following expression for the density function $f_{T\sharp\rho}$ of $T\sharp\rho$:

$$f_{T\sharp\rho}(\boldsymbol{y}) = \frac{f_{\rho}(T^{-1}(\boldsymbol{y}))}{\det\left(\nabla^{2}\varphi(T^{-1}(\boldsymbol{y}))\right)} \qquad \forall \boldsymbol{y} \in \mathbb{R}^{d}.$$
(3.1)

Since $\int_{\mathbb{R}^d} \|\boldsymbol{y}\|^2 T \sharp \rho(\mathrm{d}\boldsymbol{y}) = \int_{\mathbb{R}^d} \left\| T(\boldsymbol{x}) \right\|^2 \rho(\mathrm{d}\boldsymbol{x}) \leq \|T\|_{\mathcal{C}_{\mathrm{lin}}(\mathbb{R}^d,\mathbb{R}^d)}^2 \int_{\mathbb{R}^d} \left(1 + \|\boldsymbol{x}\|\right)^2 \rho(\mathrm{d}\boldsymbol{x}) < \infty$, we can conclude that $T \sharp \rho \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$. The proof of statement (ii) is complete.

To prove statement (iii), let us fix an arbitrary $\rho \in \mathcal{M}_{\text{full}}^q(\mathbb{R}^d)$ and denote its density function by f_{ρ} . It thus holds that $f_{\rho}(\boldsymbol{x}) > 0 \ \forall \boldsymbol{x} \in \mathbb{R}^d$, and that $f_{\rho} \in \mathcal{C}^{\text{loc},q,\alpha'}(\mathbb{R}^d)$ for some $\alpha' \in (0,1]$. By replacing α with $\min\{\alpha, \alpha'\} \in (0, 1]$ if necessary, we assume without loss of generality that $f_{\rho} \in \mathcal{C}^{\mathrm{loc}, q, \alpha}(\mathbb{R}^d)$. Since $\mathfrak{C}^{\mathrm{loc},q,\alpha}_{\underline{\lambda},\infty}(\mathbb{R}^d) \subset \mathfrak{C}^2_{\underline{\lambda},\infty}(\mathbb{R}^d)$ and $\mathcal{M}^q_{\mathrm{full}}(\mathbb{R}^d) \subset \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$, statement (ii) implies that T is a homeomorphism and $T \sharp \rho \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$, where the density function $f_{T\sharp\rho}$ of $T\sharp\rho$ is given by (3.1). Observe that (3.1) shows that $f_{T\sharp\rho}(\boldsymbol{y}) > 0$ for all $\boldsymbol{y} \in \mathbb{R}^d$. It remains to show the local Hölder property of $f_{T\sharp\rho}$. To that end, let φ^* denote the convex conjugate of φ . It follows from the duality between smooth convex functions and strongly convex functions (see, e.g., the equivalence between (a) and (e) in [56, Proposition 12.60]) and the inverse function theorem (see, e.g., [24, Theorem 1A.1]) that T^{-1} is continuously differentiable and

$$\nabla^2 \varphi^*(\boldsymbol{y}) = \nabla T^{-1}(\boldsymbol{y}) = \left[\nabla^2 \varphi \left(T^{-1}(\boldsymbol{y})\right)\right]^{-1} \qquad \forall \boldsymbol{y} \in \mathbb{R}^d.$$
(3.2)

On the one hand, since $\varphi \in \mathfrak{C}^{\mathrm{loc},q+2,\alpha}_{\underline{\lambda},\infty}(\mathbb{R}^d) \subset \mathcal{C}^{\mathrm{loc},q+2,\alpha}(\mathbb{R}^d)$, it follows from (3.2), Faà di Bruno's formula (see, e.g., [18]), and an inductive argument that $\varphi^* \in \mathcal{C}^{\mathrm{loc},q+2,\alpha}(\mathbb{R}^d)$. Consequently, since $\det(\cdot) : \mathbb{S}^d \to \mathbb{R}$ is a polynomial in all entries of the input matrix, we have by (3.2) that $\frac{1}{\det \left(\nabla^2 \varphi \left(T^{-1}(\cdot)\right)\right)} = \det \circ \nabla^2 \varphi^* \in$ $\mathcal{C}^{\mathrm{loc},q,\alpha}(\mathbb{R}^d)$. On the other hand, since $f_{\rho} \in \mathcal{C}^{\mathrm{loc},q,\alpha}(\mathbb{R}^d)$, $T^{-1} = \nabla \varphi^*$, and $\varphi^* \in \mathcal{C}^{\mathrm{loc},q+2,\alpha}(\mathbb{R}^d)$, we have by a

similar derivation using Faà di Bruno's formula and an inductive argument that $f_{\rho} \circ T^{-1} \in \mathcal{C}^{\mathrm{loc},q,\alpha}(\mathbb{R}^d)$. Hence, we conclude that $f_{T\sharp\rho} \in \mathcal{C}^{\mathrm{loc},q,\alpha}(\mathbb{R}^d)$, which completes the proof of statement (iii).

Lastly, let us prove statement (iv). Since for k = 1, ..., K, $\widehat{T}^{\mu,m_k}_{\nu_k,n_k}[\theta_k]$ satisfies the shape condition in Assumption 3.4(i), it holds \mathbb{P} -almost surely that there exist $\alpha_k \in (0,1]$, $\underline{\lambda}_k > 0$, and $\widehat{\varphi}_k \in \mathfrak{C}^{\mathrm{loc},q+2,\alpha_k}_{\underline{\lambda}_k,\infty}(\mathbb{R}^d)$ such

Algorithm 2: Computationally tractable stochastic fixed-point iterative scheme.⁵

Input: $K \in \mathbb{N}$ input probability measures $\nu_1, \ldots, \nu_K \in \mathcal{M}^q(\mathbb{R}^d)$, weights $w_1 > 0, \ldots, w_K > 0$ with $\sum_{k=1}^{K} w_k = 1$, initial probability measure $\rho_0 \in \mathcal{M}^q_{\text{full}}(\mathbb{R}^d)$, family of increasing sets $(\mathcal{X}_r)_{r \in \mathbb{N}}$, OT map estimator $\widehat{T}_{\nu,n}^{\mu,m}[\theta]$. **Output:** $(\widehat{\mu}_t)_{t \in \mathbb{N}_0}$. 1 Initialize $\hat{\rho}_0 \leftarrow \rho_0$. **2** for $t = 0, 1, 2, \dots$ do [Iteration t]: Choose $R_t \in \mathbb{N}$ using all available information up to iteration t - 1. 3 4 $\widehat{\mu}_t \leftarrow \widehat{\rho}_t |_{\mathcal{X}_{\widehat{R}_4}}.$ for $k = 1, \ldots, K$ do 5 Choose $\widehat{N}_{t,k} \in \mathbb{N}$ and $\widehat{\Theta}_{t,k} \in \Theta$ using all available information up to iteration t-1. 6 Randomly generate $\widehat{N}_{t,k}$ independent samples $\{X_{t+1,k,i}\}_{i=1:\widehat{N}_{t,k}}$ from $\widehat{\mu}_t$. 7 Randomly generate $\hat{N}_{t,k}$ independent samples $\{Y_{t+1,k,i}\}_{i=1:\hat{N}_{t,k}}$ from ν_k . 8 $\widehat{T}_{t+1,k} \leftarrow \widehat{T}_{\nu_k,\widehat{N}_{t,k}}^{\widehat{\mu}_t,\widehat{N}_{t,k}} \big[\boldsymbol{X}_{t+1,k,1}, \dots, \boldsymbol{X}_{t+1,k,\widehat{N}_{t,k}}, \boldsymbol{Y}_{t+1,k,1}, \dots, \boldsymbol{Y}_{t+1,k,\widehat{N}_{t,k}}, \widehat{\Theta}_{t,k} \big].$ 9 $\widehat{\rho}_{t+1} \leftarrow \Big[\sum_{k=1}^{K} w_k \widehat{T}_{t+1,k}\Big] \sharp \widehat{\rho}_t.$ 10 11 return $(\widehat{\mu}_t)_{t \in \mathbb{N}_0}$.

that $\widehat{T}_{\nu_k,n_k}^{\mu,m_k}[\theta_k] = \nabla \widehat{\varphi}_k$. Subsequently, let us denote $\overline{\varphi} := \sum_{k=1}^K w_k \widehat{\varphi}_k$. It follows that $\nabla \overline{\varphi} = \overline{T} \in \mathcal{C}_{\text{lin}}(\mathbb{R}^d, \mathbb{R}^d)$ and $\overline{\varphi} \in \mathfrak{C}_{\underline{\lambda},\infty}^{\text{loc},q+2,\underline{\alpha}}(\mathbb{R}^d)$ for $\underline{\alpha} := \min_{1 \le k \le K} \{\alpha_k\} \in (0,1]$ and $\underline{\lambda} := \min_{1 \le k \le K} \{\underline{\lambda}_k\} > 0$. Thus, statement (iv) follows from statement (iii). The proof is now complete.

With the above notions and properties, Algorithm 2 describes a computationally tractable algorithm which completes Conceptual Algorithm 1. The setting for Algorithm 2 is presented below.

Setting 3.6 (Inputs of Algorithm 2). In the inputs of Algorithm 2, we assume that $\nu_1, \ldots, \nu_K \in \mathcal{M}^q(\mathbb{R}^d)$ for $q \in \mathbb{N}_0$, and the weights $w_1 > 0, \ldots, w_K > 0$ satisfy $\sum_{k=1}^K w_k = 1$. ρ_0 is an arbitrary probability measure in $\mathcal{M}^q_{\text{full}}(\mathbb{R}^d)$. Moreover, we assume that $(\mathcal{X}_r)_{r\in\mathbb{N}}$ is a family of increasing sets satisfying the conditions in Assumption 3.3, and $\widehat{T}^{\mu,m}_{\nu,n}[\theta]$ is an OT map estimator satisfying the conditions in Assumption 3.4. Furthermore, we assume that $\widehat{N}_{t,k} \ge \max{\underline{m},\underline{n}} \ \forall 1 \le k \le K, \ \forall t \in \mathbb{N}_0 \ \mathbb{P}$ -almost surely, where $\underline{m},\underline{n} \in \mathbb{N}$ are given by Assumption 3.4.

Let us examine the stochastic processes generated by Algorithm 2. To begin, let us consider a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ on which the random samples in Line 7 and Line 8 are defined. Let $\mathcal{F}_0 := \{\emptyset, \Omega\}$. Observe that $\hat{\rho}_0 : \Omega \to \mathcal{M}^q_{\text{full}}(\mathbb{R}^d)$ initialized in Line 1 takes a pre-specified value ρ_0 and is thus \mathcal{F}_0 -measurable. We iteratively generate a filtration $(\mathcal{F}_t)_{t\in\mathbb{N}_0}$ as follows. Suppose that $\hat{\rho}_t$ is \mathcal{F}_t -measurable for some $t \in \mathbb{N}_0$. We let the index $\hat{R}_t : \Omega \to \mathbb{N}$ in Line 3 be an \mathcal{F}_t -measurable random variable. After \hat{R}_t has been chosen, Proposition 3.5(i) implies that $\hat{\mu}_t := \hat{\rho}_t|_{\mathcal{X}_{\hat{R}_t}}$ in Line 4 is an $\mathcal{M}^q(\mathbb{R}^d)$ -valued random variable, which is also \mathcal{F}_t -measurable. Subsequently, for $k = 1, \ldots, K$, we let the sample size $\hat{N}_{t,k} : \Omega \to \mathbb{N}$ and the parameter $\hat{\Theta}_{t,k} : \Omega \to \Theta$ in Line 6 be \mathcal{F}_t -measurable random variables. After $\hat{N}_{t,k}$ and $\hat{\Theta}_{t,k}$ have been chosen, $\hat{N}_{t,k}$ independent samples $\mathbf{X}_{t+1,k,1}, \ldots, \mathbf{X}_{t+1,k,\hat{N}_{t,k}} : \Omega \to \mathbb{R}^d$ from $\hat{\mu}_t$ and $\hat{N}_{t,k}$ independent samples $\mathbf{X}_{t+1,k,1}, \ldots, \mathbf{X}_{t+1,k,\hat{N}_{t,k}} : \Omega \to \mathbb{R}^d$ from $\hat{\mu}_t$ and $\hat{N}_{t,k}$ independent samples $\{\mathbf{X}_{t+1,k,1}, \ldots, \mathbf{X}_{t+1,k,\hat{N}_{t,k}}\}_{k=1:K}$ to be jointly independent conditional on \mathcal{F}_t . Let

⁵The input configuration of this algorithm is specified in Setting 3.6. Concrete choices of the truncation index (Line 3), the sample size (Line 6), and the parameter of the OT map estimator (Line 6) to ensure convergence of the output sequence of probability measures are specified in Setting 3.13.

 \mathcal{F}_{t+1} be the σ -algebra generated by all the random samples up to iteration t, i.e.,

$$\mathcal{F}_{t+1} := \sigma \left(\bigcup_{0 \le s \le t} \left(\{ \boldsymbol{X}_{s+1,k,i} \}_{i=1:\widehat{N}_{s,k}, \, k=1:K} \cup \{ \boldsymbol{Y}_{s+1,k,i} \}_{i=1:\widehat{N}_{s,k}, \, k=1:K} \right) \right) \qquad \forall t \in \mathbb{N}_0.$$

For k = 1, ..., K, the OT map estimator $\widehat{T}_{t+1,k} : \Omega \to C_{\text{lin}}(\mathbb{R}^d, \mathbb{R}^d)$ in Line 9 is thus \mathcal{F}_{t+1} -measurable. Since $\widehat{N}_{t,k} \geq \max\{\underline{m},\underline{n}\}$ for k = 1, ..., K P-almost surely, Proposition 3.5(iv) guarantees that $\widehat{\rho}_{t+1} := \left[\sum_{k=1}^{K} w_k \widehat{T}_{t+1,k}\right] \sharp \widehat{\rho}_t$ in Line 10 is an $\mathcal{M}_{\text{full}}^q(\mathbb{R}^d)$ -valued random variable. Since $\widehat{\rho}_{t+1}$ depends on $\widehat{\rho}_t$ and $(\widehat{T}_{t+1,k})_{k=1:K}$, it is \mathcal{F}_{t+1} -measurable. Iteratively repeating the above construction for t = 0, 1, 2, ... leads to a filtered probability space with filtration $(\mathcal{F}_t)_{t\in\mathbb{N}_0}$. The resulting sequences $(\widehat{\rho}_t)_{t\in\mathbb{N}_0}$ and $(\widehat{\mu}_t)_{t\in\mathbb{N}_0}$ are thus $(\mathcal{F}_t)_{t\in\mathbb{N}_0}$ -adapted stochastic processes. In the next subsection, we will specify the choices of $(\widehat{R}_t)_{t\in\mathbb{N}_0}$, $(\widehat{N}_{t,k})_{k=1:K,t\in\mathbb{N}_0}$, and $(\widehat{\Theta}_{t,k})_{k=1:K,t\in\mathbb{N}_0}$ (which are required to be $(\mathcal{F}_t)_{t\in\mathbb{N}_0}$ -adapted stochastic processes) in order to achieve P-almost sure convergence of the output process $(\widehat{\mu}_t)_{t\in\mathbb{N}_0}$ of Algorithm 2.

Remark 3.7. In Algorithm 2, rather than directly updating $\hat{\mu}_t$ to $\hat{\mu}_{t+1} \leftarrow \left[\sum_{k=1}^K w_k \hat{T}_{t+1,k}\right] \sharp \hat{\mu}_t$, we first apply the pushforward of $\hat{\rho}_t \in \mathcal{M}^q_{\text{full}}(\mathbb{R}^d)$ by $\left[\sum_{k=1}^K w_k \hat{T}_{t+1,k}\right]$ in Line 10 to obtain $\hat{\rho}_{t+1} \in \mathcal{M}^q_{\text{full}}(\mathbb{R}^d)$, and then truncate $\hat{\rho}_{t+1}$ to $\mathcal{X}_{\hat{R}_{t+1}}$ to get $\hat{\mu}_{t+1}$. The truncation step guarantees that $\hat{\mu}_{t+1} \in \mathcal{M}^q(\mathbb{R}^d)$ so that the consistency condition of the OT map estimator in Assumption 3.4(iii) can be satisfied (see our results and discussions in Section 4). Note that the support of the pushforward $\left[\sum_{k=1}^K w_k \hat{T}_{t+1,k}\right] \sharp \hat{\mu}_t$ is not necessarily an admissible support set in $\mathcal{S}(\mathbb{R}^d)$; specifically, the uniform convexity condition may fail.

Remark 3.8 (Computational tractability of Algorithm 2). We assume that: (i) independent random samples from ν_1, \ldots, ν_K , and ρ_0 can be efficiently generated; (ii) the OT map estimator $\widehat{T}_{\nu,n}^{\mu,m}[\theta]$ can be tractably computed and $\widehat{T}_{\nu,n}^{\mu,m}[\theta](\mathbf{x})$ can be tractably evaluated at any point $\mathbf{x} \in \mathbb{R}^d$; (iii) for all $r \in \mathbb{N}$, checking whether a point $\mathbf{x} \in \mathbb{R}^d$ belongs to \mathcal{X}_r is computationally tractable. Then, Algorithm 2 is computationally tractable. Indeed, for $t \in \mathbb{N}$, a random sample from $\widehat{\mu}_t$ can be generated by rejection sampling. Specifically, one first generates a random sample $\mathbf{X} \in \mathbb{R}^d$ from ρ_0 and evaluates the composition $\widehat{\mathbf{X}} := \left[\sum_{k=1}^K w_k \widehat{T}_{t,k}\right] \circ$ $\cdots \circ \left[\sum_{k=1}^K w_k \widehat{T}_{1,k}\right](\mathbf{X})$. This sample $\widehat{\mathbf{X}}$ is subsequently accepted if $\widehat{\mathbf{X}} \in \mathcal{X}_{\widehat{R}_t}$. Otherwise, it is generated repeatedly until accepted. The computational tractability of our specific choice of OT map estimator is discussed in Remark 4.4 in Section 4.

Remark 3.9 (Distributed implementation of Algorithm 2). We would like to remark that our proposed Algorithm 2 allows for implementation in a distributed and parallel computing environment, which can be appealing in terms of computational efficiency. Suppose that there are a large number K of agents each having local access to an input measure $\nu_k \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$. The \mathcal{W}_2 -barycenter problem instance with input measures ν_1, \ldots, ν_K and weights w_1, \ldots, w_K is to be solved by a central coordinator who can communicate with the K agents. In each iteration t, the coordinator first generates independent samples $\{\mathbf{X}_{t+1,k,i}\}_{i=1:\widehat{N}_{t,k}}$ from $\hat{\mu}_t$ and release the subcollection of samples $\{\mathbf{X}_{t+1,k,i}\}_{i=1:\widehat{N}_{t,k}}$ from ν_k (Line 7), for $k = 1, \ldots, K$. Each agent k then generates independent samples $\{\mathbf{Y}_{t+1,k,i}\}_{i=1:\widehat{N}_{t,k}}$ from ν_k (Line 8) and uses $\{\mathbf{X}_{t+1,k,i}\}_{i=1:\widehat{N}_{t,k}}$ and $\{\mathbf{Y}_{t+1,k,i}\}_{i=1:\widehat{N}_{t,k}}$ to compute an admissible OT map estimator $\widehat{T}_{t+1,k}$ (Line 9). Subsequently, in order for the coordinator to generate independent samples from $\widehat{\mu}_{t+1}$ conditional on \mathcal{F}_t , a large number N of independent samples $\{\widehat{X}_{t+1,i}\}_{i=1:N}$ from the coordinator, agent k evaluates $\{\widehat{T}_{t+1,k}(\widehat{X}_{t+1,i})\}_{i=1:N}$ and sends it back to the coordinator. The coordinator can then generate independent samples from $\widehat{\mu}_{t+1}$ using the weighted sums: $\{\sum_{k=1}^K w_k \widehat{T}_{t+1,k}(\widehat{X}_{t+1,i})\}_{i=1:N}$ (Line 10) followed by the rejection sampling procedure described in Remark 3.8.

3.2. Convergence analysis. The goal of this subsection is to develop sufficient conditions for the convergence of the output process $(\hat{\mu}_t)_{t\in\mathbb{N}_0}$ in Algorithm 2. Let us begin by analyzing the decrements of the process $(V(\hat{\mu}_t))_{t\in\mathbb{N}_0}$. This will subsequently lead to sufficient conditions on the choices of $(\hat{R}_t)_{t\in\mathbb{N}_0}, (\hat{N}_{t,k})_{k=1:K, t\in\mathbb{N}_0}$, and $(\hat{\Theta}_{t,k})_{k=1:K, t\in\mathbb{N}_0}$ to guarantee the \mathbb{P} -almost sure convergence of $(\hat{\mu}_t)_{t\in\mathbb{N}_0}$. **Proposition 3.10** (Decrement of the process $(V(\hat{\mu}_t))_{t \in \mathbb{N}_0}$). Let the inputs of Algorithm 2 satisfy Setting 3.6, let $(\Omega, \mathcal{F}, \mathbb{P}, (\mathcal{F}_t)_{t \in \mathbb{N}_0})$ be the filtered probability space generated by Algorithm 2, and let $(\hat{\mu}_t)_{t \in \mathbb{N}_0}$ be the output of Algorithm 2. Moreover, let $V(\cdot)$ be the function defined in (1.2) and let $G(\cdot)$ be the operator defined in (1.3). Then, the sequence $(V(\hat{\mu}_t))_{t \in \mathbb{N}_0}$ satisfies

$$V(\widehat{\mu}_{t+1}) - V(\widehat{\mu}_{t}) \leq -\mathcal{W}_{2}(\widehat{\mu}_{t}, G(\widehat{\mu}_{t}))^{2} + 2\sum_{k=1}^{K} w_{k} \|\widehat{T}_{t+1,k} - T_{\nu_{k}}^{\widehat{\mu}_{t}}\|_{\mathcal{L}^{2}(\widehat{\mu}_{t})}^{2} + 2\mathcal{W}_{2}(\left[\sum_{k=1}^{K} w_{k}\widehat{T}_{t+1,k}\right] \sharp \widehat{\mu}_{t}, \widehat{\mu}_{t+1})^{2} \quad \forall t \in \mathbb{N}_{0}, \ \mathbb{P}\text{-}a.s.$$
(3.3)

In particular, taking conditional expectations with respect to \mathcal{F}_t on both sides of (3.3) yields

$$\mathbb{E}\left[V(\widehat{\mu}_{t+1})\big|\mathcal{F}_{t}\right] - V(\widehat{\mu}_{t}) \leq -\mathcal{W}_{2}\left(\widehat{\mu}_{t}, G(\widehat{\mu}_{t})\right)^{2} + 2\sum_{k=1}^{K} w_{k}\mathbb{E}\left[\left\|\widehat{T}_{t+1,k} - T_{\nu_{k}}^{\widehat{\mu}_{t}}\right\|_{\mathcal{L}^{2}(\widehat{\mu}_{t})}^{2}\Big|\mathcal{F}_{t}\right] + 2\mathbb{E}\left[\mathcal{W}_{2}\left(\left[\sum_{k=1}^{K} w_{k}\widehat{T}_{t+1,k}\right]\sharp\widehat{\mu}_{t}, \widehat{\mu}_{t+1}\right)^{2}\Big|\mathcal{F}_{t}\right] \quad \forall t \in \mathbb{N}_{0}, \ \mathbb{P}\text{-}a.s.$$
(3.4)

Proof of Proposition 3.10. Throughout this proof, let us fix an arbitrary $t \in \mathbb{N}_0$, denote $\overline{T}^{\widehat{\mu}_t} := \sum_{k=1}^K w_k T^{\widehat{\mu}_k}_{\nu_k}$, $\overline{T}_{t+1} := \sum_{k=1}^K w_k \widehat{T}_{t+1,k}$, and denote $\widetilde{\mu}_{t+1} := \overline{T}_{t+1} \sharp \widehat{\mu}_t$. Since Setting 3.6 guarantees $\widehat{N}_{t,k} \ge \min\{\underline{m},\underline{n}\}$ for $k = 1, \ldots, K$ P-almost surely, we have by Proposition 3.5(iv) that $\overline{T}_{t+1} \in \mathcal{C}_{\text{lin}}(\mathbb{R}^d, \mathbb{R}^d)$ is P-almost surely a homeomorphism. Subsequently, Proposition 3.5(ii) implies that $\widetilde{\mu}_{t+1} \in \mathcal{P}_{2,\text{ac}}(\mathbb{R}^d)$ P-almost surely. Let $T^{\widetilde{\mu}_{t+1}}_{\widehat{\mu}_{t+1}} : \mathbb{R}^d \to \mathbb{R}^d$ denotes the OT map from $\widetilde{\mu}_{t+1}$ to $\widehat{\mu}_{t+1}$ which exists and is $\widetilde{\mu}_{t+1}$ -almost everywhere unique due to Brenier's theorem (Theorem 2.4). In the remainder of this proof, all statements hold in the P-almost sure sense, and we will omit "P-a.s." for ease of notation.

Our proof below uses the following identity, which can be verified directly by expanding both sides:

$$\sum_{k=1}^{K} w_k \| \boldsymbol{y} - \boldsymbol{z}_k \|^2 = \| \boldsymbol{y} - \bar{\boldsymbol{z}} \|^2 + \sum_{k=1}^{K} w_k \| \bar{\boldsymbol{z}} - \boldsymbol{z}_k \|^2$$

$$\text{where } \bar{\boldsymbol{z}} := \sum_{k=1}^{K} w_k \boldsymbol{z}_k \qquad \forall \boldsymbol{y}, \boldsymbol{z}_1, \dots, \boldsymbol{z}_k \in \mathbb{R}^d.$$
(3.5)

For any $x \in \mathbb{R}^d$, substituting $y \leftarrow x$ and $z_k \leftarrow T_{\nu_k}^{\widehat{\mu}_t}(x)$ in (3.5) gives us

$$\sum_{k=1}^{K} w_k \left\| \boldsymbol{x} - T_{\nu_k}^{\hat{\mu}_t}(\boldsymbol{x}) \right\|^2 = \left\| \boldsymbol{x} - \bar{T}^{\hat{\mu}_t}(\boldsymbol{x}) \right\|^2 + \sum_{k=1}^{K} w_k \left\| \bar{T}^{\hat{\mu}_t}(\boldsymbol{x}) - T_{\nu_k}^{\hat{\mu}_t}(\boldsymbol{x}) \right\|^2.$$
(3.6)

Moreover, substituting $\boldsymbol{y} \leftarrow T_{\hat{\mu}_{t+1}}^{\tilde{\mu}_{t+1}} \circ \bar{T}_{t+1}(\boldsymbol{x})$ and $\boldsymbol{z}_k \leftarrow T_{\nu_k}^{\hat{\mu}_t}(\boldsymbol{x})$ in (3.5), we obtain

$$\sum_{k=1}^{K} w_k \left\| T_{\hat{\mu}_{t+1}}^{\tilde{\mu}_{t+1}} \circ \bar{T}_{t+1}(\boldsymbol{x}) - T_{\nu_k}^{\hat{\mu}_t}(\boldsymbol{x}) \right\|^2 = \left\| T_{\hat{\mu}_{t+1}}^{\tilde{\mu}_{t+1}} \circ \bar{T}_{t+1}(\boldsymbol{x}) - \bar{T}^{\hat{\mu}_t}(\boldsymbol{x}) \right\|^2 + \sum_{k=1}^{K} w_k \left\| \bar{T}^{\hat{\mu}_t}(\boldsymbol{x}) - T_{\nu_k}^{\hat{\mu}_t}(\boldsymbol{x}) \right\|^2.$$

Combining this with (3.6) yields

$$\left(\sum_{k=1}^{K} w_{k} \left\| T_{\widehat{\mu}_{t+1}}^{\widetilde{\mu}_{t+1}} \circ \bar{T}_{t+1}(\boldsymbol{x}) - T_{\nu_{k}}^{\widehat{\mu}_{t}}(\boldsymbol{x}) \right\|^{2} \right) - \left(\sum_{k=1}^{K} w_{k} \left\| \boldsymbol{x} - T_{\nu_{k}}^{\widehat{\mu}_{t}}(\boldsymbol{x}) \right\|^{2} \right) \\
= \left\| T_{\widehat{\mu}_{t+1}}^{\widetilde{\mu}_{t+1}} \circ \bar{T}_{t+1}(\boldsymbol{x}) - \bar{T}^{\widehat{\mu}_{t}}(\boldsymbol{x}) \right\|^{2} - \left\| \boldsymbol{x} - \bar{T}^{\widehat{\mu}_{t}}(\boldsymbol{x}) \right\|^{2} \\
\leq 2 \left\| T_{\widehat{\mu}_{t+1}}^{\widetilde{\mu}_{t+1}} \circ \bar{T}_{t+1}(\boldsymbol{x}) - \bar{T}_{t+1}(\boldsymbol{x}) \right\|^{2} + 2 \left\| \bar{T}_{t+1}(\boldsymbol{x}) - \bar{T}^{\widehat{\mu}_{t}}(\boldsymbol{x}) \right\|^{2} - \left\| \boldsymbol{x} - \bar{T}^{\widehat{\mu}_{t}}(\boldsymbol{x}) \right\|^{2} - \left\| \boldsymbol{x} - \bar{T}^{\widehat{\mu}_{t}}(\boldsymbol{x}) \right\|^{2} \quad \forall \boldsymbol{x} \in \mathbb{R}^{d}.$$
(3.7)

In the following, let us examine the integral of each term in (3.7) with respect to $\hat{\mu}_t$. Firstly, for $k = 1, \ldots, K$, let $\pi_k := \left[T_{\hat{\mu}_{t+1}}^{\tilde{\mu}_{t+1}} \circ \bar{T}_{t+1}, T_{\nu_k}^{\hat{\mu}_t} \right] \sharp \hat{\mu}_t \in \mathcal{P}(\mathbb{R}^d \times \mathbb{R}^d)$. Since $\left(T_{\hat{\mu}_{t+1}}^{\tilde{\mu}_{t+1}} \circ \bar{T}_{t+1} \right) \sharp \hat{\mu}_t = T_{\hat{\mu}_{t+1}}^{\tilde{\mu}_{t+1}} \sharp \tilde{\mu}_{t+1} = \hat{\mu}_{t+1}$ and $T_{\nu_k}^{\widehat{\mu}_t} \sharp \widehat{\mu}_t = \nu_k$, it follows that $\pi_k \in \Pi(\widehat{\mu}_{t+1}, \nu_k)$. Thus, we get

$$\sum_{k=1}^{K} w_k \int_{\mathbb{R}^d} \left\| T_{\hat{\mu}_{t+1}}^{\tilde{\mu}_{t+1}} \circ \bar{T}_{t+1}(\boldsymbol{x}) - T_{\nu_k}^{\hat{\mu}_t}(\boldsymbol{x}) \right\|^2 \hat{\mu}_t(\mathrm{d}\boldsymbol{x}) = \sum_{k=1}^{K} w_k \int_{\mathbb{R}^d \times \mathbb{R}^d} \|\boldsymbol{x}_1 - \boldsymbol{x}_2\|^2 \pi_k(\mathrm{d}\boldsymbol{x}_1, \mathrm{d}\boldsymbol{x}_2)$$

$$\geq \sum_{k=1}^{K} w_k \mathcal{W}_2(\hat{\mu}_{t+1}, \nu_k)^2 = V(\hat{\mu}_{t+1}).$$
(3.8)

Secondly, since $T_{\nu_k}^{\hat{\mu}_t}$ is the OT map from $\hat{\mu}_t$ to ν_k for $k = 1, \ldots, K$, we have

$$\sum_{k=1}^{K} w_k \int_{\mathbb{R}^d} \left\| \boldsymbol{x} - T_{\nu_k}^{\widehat{\mu}_t}(\boldsymbol{x}) \right\|^2 \widehat{\mu}_t(\mathrm{d}\boldsymbol{x}) = \sum_{k=1}^{K} w_k \mathcal{W}_2(\widehat{\mu}_t, \nu_k)^2 = V(\widehat{\mu}_t).$$
(3.9)

Thirdly, since $T_{\widehat{\mu}_{t+1}}^{\widetilde{\mu}_{t+1}}$ is the OT map from $\widetilde{\mu}_{t+1}$ to $\widehat{\mu}_{t+1}$, it holds that

$$\int_{\mathbb{R}^d} \left\| T_{\hat{\mu}_{t+1}}^{\tilde{\mu}_{t+1}} \circ \bar{T}_{t+1}(\boldsymbol{x}) - \bar{T}_{t+1}(\boldsymbol{x}) \right\|^2 \hat{\mu}_t(\mathrm{d}\boldsymbol{x}) = \int_{\mathbb{R}^d} \left\| T_{\hat{\mu}_{t+1}}^{\tilde{\mu}_{t+1}}(\boldsymbol{y}) - \boldsymbol{y} \right\|^2 \tilde{\mu}_{t+1}(\mathrm{d}\boldsymbol{y}) = \mathcal{W}_2(\tilde{\mu}_{t+1}, \hat{\mu}_{t+1})^2.$$
(3.10)

Fourthly, the convexity of $\mathbb{R}^d \ni z \mapsto ||z||^2 \in \mathbb{R}$ together with Jensen's inequality gives

$$\left\|\bar{T}_{t+1}(\boldsymbol{x}) - \bar{T}^{\widehat{\mu}_{t}}(\boldsymbol{x})\right\|^{2} = \left\|\sum_{k=1}^{K} w_{k} \left(\widehat{T}_{t+1,k}(\boldsymbol{x}) - T^{\widehat{\mu}_{t}}_{\nu_{k}}(\boldsymbol{x})\right)\right\|^{2} \le \sum_{k=1}^{K} w_{k} \left\|\widehat{T}_{t+1,k}(\boldsymbol{x}) - T^{\widehat{\mu}_{t}}_{\nu_{k}}(\boldsymbol{x})\right\|^{2} \quad \forall \boldsymbol{x} \in \mathbb{R}^{d},$$

which results in

$$\int_{\mathbb{R}^{d}} \left\| \bar{T}_{t+1}(\boldsymbol{x}) - \bar{T}^{\widehat{\mu}_{t}}(\boldsymbol{x}) \right\|^{2} \widehat{\mu}_{t}(\mathrm{d}\boldsymbol{x}) \leq \sum_{k=1}^{K} w_{k} \int_{\mathbb{R}^{d}} \left\| \widehat{T}_{t+1,k}(\boldsymbol{x}) - T^{\widehat{\mu}_{t}}_{\nu_{k}}(\boldsymbol{x}) \right\|^{2} \widehat{\mu}_{t}(\mathrm{d}\boldsymbol{x})
= \sum_{k=1}^{K} w_{k} \left\| \widehat{T}_{t+1,k} - T^{\widehat{\mu}_{t}}_{\nu_{k}} \right\|_{\mathcal{L}^{2}(\widehat{\mu}_{t})}^{2}.$$
(3.11)

Lastly, for k = 1, ..., K, let $\varphi_{\nu_k}^{\hat{\mu}_t}$ denote the Brenier potential from $\hat{\mu}_t$ to ν_k , which is a proper, l.s.c., and convex function on \mathbb{R}^d . Since $\bar{T}^{\hat{\mu}_t}$ is $\hat{\mu}_t$ -almost everywhere equal to the gradient of the proper, l.s.c., and convex function $\sum_{k=1}^K w_k \varphi_{\nu_k}^{\hat{\mu}_t}$, it follows from Brenier's theorem (Theorem 2.4) that $\bar{T}^{\hat{\mu}_t}$ is the OT map from $\hat{\mu}_t$ to $\bar{T}^{\hat{\mu}_t} \sharp \hat{\mu}_t = G(\hat{\mu}_t)$, resulting in

$$\int_{\mathbb{R}^d} \left\| \boldsymbol{x} - \bar{T}^{\widehat{\mu}_t}(\boldsymbol{x}) \right\|^2 \widehat{\mu}_t(\mathrm{d}\boldsymbol{x}) = \mathcal{W}_2(\widehat{\mu}_t, G(\widehat{\mu}_t))^2.$$
(3.12)

Now, integrating both sides of (3.7) with respect to $\hat{\mu}_t$ and then combining it with (3.8)–(3.12) completes the proof of (3.3). Finally, taking conditional expectations with respect to \mathcal{F}_t on both sides of (3.3) proves (3.4). The proof is now complete.

Remark 3.11. In [3, Proposition 3.3], the decrement of the sequence $(V(\mu_t))_{t\in\mathbb{N}_0}$ in the deterministic fixedpoint iteration $\mu_{t+1} \leftarrow G(\mu_t) \ \forall t \in \mathbb{N}_0$ is controlled through the inequality:

$$V(\mu_{t+1}) - V(\mu_t) \le -\mathcal{W}_2(\mu_t, G(\mu_t))^2 \qquad \forall t \in \mathbb{N}_0.$$
(3.13)

Compared to (3.13), the stochastic decrement (3.3) in Proposition 3.10 has two additional terms on the righthand side:

- the term $2\sum_{k=1}^{K} w_k \|\widehat{T}_{t+1,k} T_{\nu_k}^{\widehat{\mu}_t}\|_{\mathcal{L}^2(\widehat{\mu}_t)}^2$ comes from the inexactness when approximating the true OT map $T_{\nu_k}^{\widehat{\mu}_t}$ by the OT map estimator $\widehat{T}_{t+1,k}$, i.e., from the approximation in Line 6 of Conceptual Algorithm 1;
- the term $2W_2(\left[\sum_{k=1}^K w_k \widehat{T}_{t+1,k}\right] \sharp \widehat{\mu}_t, \widehat{\mu}_{t+1})^2$ comes from the inexactness when approximating the pushforward $\left[\sum_{k=1}^K w_k \widehat{T}_{t+1,k}\right] \sharp \widehat{\mu}_t$ by $\widehat{\mu}_{t+1} = \left(\left[\sum_{k=1}^K w_k \widehat{T}_{t+1,k}\right] \sharp \widehat{\rho}_t\right)\Big|_{\mathcal{X}_{\widehat{R}_{t+1}}}$, i.e., from the approximation

in Line 7 of Conceptual Algorithm 1.

In order to guarantee the convergence of the output process $(\hat{\mu}_t)_{t\in\mathbb{N}_0}$ of Algorithm 2, we aim to control the two error terms $2\sum_{k=1}^{K} w_k \mathbb{E}\left[\|\widehat{T}_{t+1,k} - T_{\nu_k}^{\widehat{\mu}_t}\|_{\mathcal{L}^2(\widehat{\mu}_t)}^2 |\mathcal{F}_t \right]$ and $2\mathbb{E}\left[\mathcal{W}_2\left(\left[\sum_{k=1}^{K} w_k \widehat{T}_{t+1,k} \right] \sharp \widehat{\mu}_t, \widehat{\mu}_{t+1} \right)^2 |\mathcal{F}_t \right]$ on the right-hand side of (3.4) to be arbitrarily close to 0. Before presenting our concrete setting of Algorithm 2 that guarantees the convergence of $(\widehat{\mu}_t)_{t\in\mathbb{N}_0}$, let us first establish an intermediate result about choosing the truncation set $\mathcal{X}_{\widehat{R}_t}$ in Line 3 presented in the lemma below.

Lemma 3.12 (Choice of the truncation set). Let $q \in \mathbb{N}_0$, let $\nu_1, \ldots, \nu_K \in \mathcal{M}^q(\mathbb{R}^d)$, and let $\rho \in \mathcal{M}^q_{\text{full}}(\mathbb{R}^d)$. Moreover, let $(\mathcal{X}_r)_{r \in \mathbb{N}}$ be a family of increasing sets satisfying Assumption 3.3 and let $\widehat{T}^{\mu,m}_{\nu,n}[\theta]$ be an OT map estimator satisfying Assumption 3.4. Then, the following statements hold.

- (i) There exists $\overline{r}_1(\rho, \epsilon) \in \mathbb{N}$ that depends on ρ, ϵ such that for all $\epsilon > 0$ and all $r \ge \overline{r}_1(\rho, \epsilon)$, the truncated probability measure $\mu := \rho|_{\mathcal{X}_r}$ satisfies $\mathcal{W}_2(\mu, \rho)^2 \le \epsilon$.
- (ii) There exists $\overline{r}_2(\rho, \nu_1, \dots, \nu_K, \epsilon) \in \mathbb{N}$ that depends on $\rho, \nu_1, \dots, \nu_K, \epsilon$ such that for all $\epsilon > 0$ and all $r \geq \overline{r}_2(\rho, \nu_1, \dots, \nu_K, \epsilon)$, the truncated probability measure $\mu := \rho|_{\mathcal{X}_r}$ satisfies

$$\mathbb{E}\Big[\mathcal{W}_2\big(\widehat{T}^{\mu,m_k}_{\nu_k,n_k}[\theta_k] \sharp \mu, \widehat{T}^{\mu,m_k}_{\nu_k,n_k}[\theta_k] \sharp \rho\big)^2\Big] \le \epsilon \qquad \forall m_k \ge \underline{m}, \ \forall n_k \ge \underline{n}, \ \forall \theta_k \in \Theta, \ \forall 1 \le k \le K,$$

where $\underline{m}, \underline{n} \in \mathbb{N}$ are constants given by Assumption 3.4(i). Moreover, in this case, whenever $m_k \geq \underline{m}$, $n_k \geq \underline{n}, \theta_k \in \Theta$ for $k = 1, \ldots, K$, $\overline{T} := \sum_{k=1}^{K} w_k \widehat{T}_{\nu_k, n_k}^{\mu, m_k}[\theta_k]$ satisfies $\mathbb{E}\left[\mathcal{W}_2(\overline{T} \sharp \mu, \overline{T} \sharp \rho)^2 \right] \leq \epsilon$.

Proof of Lemma 3.12. Let us first prove statement (i). For every $r \in \mathbb{N}$, let us define $\hat{\mu}_r := \rho|_{\mathcal{X}_r}$ and define $\check{\mu}_r := \rho|_{\mathcal{X}_r^c}$, where $\mathcal{X}_r^c := \mathbb{R}^d \setminus \mathcal{X}_r$. Notice that $\rho = \rho(\mathcal{X}_r)\hat{\mu}_r + (1 - \rho(\mathcal{X}_r))\check{\mu}_r$ for all $r \in \mathbb{N}$. Let $\pi_{r,1} := [I_d, I_d] \sharp \hat{\mu}_r$ where $I_d : \mathbb{R}^d \to \mathbb{R}^d$ denotes the identity mapping on \mathbb{R}^d , let $\pi_{r,2} \in \Pi(\hat{\mu}_r, \check{\mu}_r)$ be arbitrary, and let $\pi_r := \rho(\mathcal{X}_r)\pi_{r,1} + (1 - \rho(\mathcal{X}_r))\pi_{r,2} \in \mathcal{P}(\mathbb{R}^d \times \mathbb{R}^d)$. One may check that $\pi_r \in \Pi(\hat{\mu}_r, \rho)$ for all $r \in \mathbb{N}$. Subsequently, it holds for all $r \in \mathbb{N}$ that

$$\begin{split} \mathcal{W}_{2}(\hat{\mu}_{r},\rho)^{2} &\leq \int_{\mathbb{R}^{d}\times\mathbb{R}^{d}} \|\boldsymbol{x}-\boldsymbol{y}\|^{2} \,\pi_{r}(\mathrm{d}\boldsymbol{x},\mathrm{d}\boldsymbol{y}) \\ &= \rho(\mathcal{X}_{r}) \int_{\mathbb{R}^{d}} \|\boldsymbol{x}-\boldsymbol{x}\|^{2} \,\hat{\mu}_{r}(\mathrm{d}\boldsymbol{x}) + (1-\rho(\mathcal{X}_{r})) \int_{\mathbb{R}^{d}\times\mathbb{R}^{d}} \|\boldsymbol{x}-\boldsymbol{y}\|^{2} \,\pi_{r,2}(\mathrm{d}\boldsymbol{x},\mathrm{d}\boldsymbol{y}) \\ &\leq (1-\rho(\mathcal{X}_{r})) \int_{\mathbb{R}^{d}\times\mathbb{R}^{d}} 2\|\boldsymbol{x}\|^{2} + 2\|\boldsymbol{y}\|^{2} \,\pi_{r,2}(\mathrm{d}\boldsymbol{x},\mathrm{d}\boldsymbol{y}) \\ &= (1-\rho(\mathcal{X}_{r})) \int_{\mathbb{R}^{d}} 2\|\boldsymbol{x}\|^{2} \,\hat{\mu}_{r}(\mathrm{d}\boldsymbol{x}) + (1-\rho(\mathcal{X}_{r})) \int_{\mathbb{R}^{d}} 2\|\boldsymbol{y}\|^{2} \,\check{\mu}_{r}(\mathrm{d}\boldsymbol{y}) \\ &\leq \frac{1-\rho(\mathcal{X}_{r})}{\rho(\mathcal{X}_{r})} \int_{\mathbb{R}^{d}} 2\|\boldsymbol{x}\|^{2} \,\rho(\mathrm{d}\boldsymbol{x}) + \int_{\mathbb{R}^{d}} 2\|\boldsymbol{y}\|^{2} \mathbb{1}_{\mathcal{X}_{r}^{c}}(\boldsymbol{y}) \,\rho(\mathrm{d}\boldsymbol{y}). \end{split}$$

Since $\bigcup_{r \in \mathbb{N}_0} \mathcal{X}_r = \mathbb{R}^d$ and $\rho \in \mathcal{P}_2(\mathbb{R}^d)$ by assumption, it follows from Lebesgue's dominated convergence theorem that

$$\limsup_{r \to \infty} \mathcal{W}_2(\hat{\mu}_r, \rho)^2 \leq \limsup_{r \to \infty} \frac{1 - \rho(\mathcal{X}_r)}{\rho(\mathcal{X}_r)} \int_{\mathbb{R}^d} 2 \|\boldsymbol{x}\|^2 \,\rho(\mathrm{d}\boldsymbol{x}) + \limsup_{r \to \infty} \int_{\mathbb{R}^d} 2 \|\boldsymbol{y}\|^2 \mathbb{1}_{\mathcal{X}_r^c}(\boldsymbol{y}) \,\rho(\mathrm{d}\boldsymbol{y}) = 0.$$

Therefore, for any $\epsilon > 0$, there exists $\overline{r}_1(\rho, \epsilon) \in \mathbb{N}$ such that $\mathcal{W}_2(\hat{\mu}_r, \rho)^2 \leq \epsilon$ for all $r \geq \overline{r}_1(\rho, \epsilon)$. This proves statement (i).

Next, let us use the growth condition in Assumption 3.4(ii) to prove statement (ii). For every $r \in \mathbb{N}$, let $\hat{\mu}_r, \check{\mu}_r, \pi_{r,1}, \pi_{r,2}$, and π_r be defined as in the proof of statement (i). Recall that $\pi_r \in \Pi(\hat{\mu}_r, \rho)$. Moreover, for $k = 1, \ldots, K$, let $m_k \geq \underline{m}, n_k \geq \underline{n}, \theta_k \in \Theta$ be arbitrary and denote $\dot{T}_{k,r} := \hat{T}_{\nu_k, n_k}^{\hat{\mu}_r, m_k}[\theta_k], \dot{T}_r := \sum_{k=1}^K w_k \dot{T}_{k,r}$ for notational simplicity. Furthermore, for $k = 1, \ldots, K$, we denote by $\dot{T}_{k,r} \otimes \dot{T}_{k,r}$ the function $\mathbb{R}^d \times \mathbb{R}^d \ni (x, y) \mapsto \dot{T}_{k,r} \otimes \dot{T}_{k,r}(x, y) := (\dot{T}_{k,r}(x), \dot{T}_{k,r}(y)) \in \mathbb{R}^d \times \mathbb{R}^d$. Similarly, we denote by $\dot{T}_r \otimes \dot{T}_r$ the function $\mathbb{R}^d \times \mathbb{R}^d \ni (x, y) \mapsto \dot{T}_r \otimes \dot{T}_r(x, y) := (\dot{T}_r(x), \dot{T}_r(y)) \in \mathbb{R}^d \times \mathbb{R}^d$. It holds that $[\dot{T}_{k,r} \otimes \dot{T}_{k,r}] \sharp \pi_r \in \Pi(\dot{T}_{k,r} \sharp \hat{\mu}_r, \dot{T}_{k,r} \ddagger \rho)$ for $k = 1, \ldots, K$, and $[\dot{T}_r \otimes \dot{T}_r] \sharp \pi_r \in \Pi(\dot{T}_r \sharp \hat{\mu}_r, \dot{T}_r \sharp \rho)$. Therefore, for $k = 1, \ldots, K$, we

are able to bound $\mathcal{W}_2(\dot{T}_{k,r} \sharp \hat{\mu}_r, \dot{T}_{k,r} \sharp \rho)^2$ by

$$\begin{aligned} \mathcal{W}_{2}(\dot{T}_{k,r}\sharp\hat{\mu}_{r},\dot{T}_{k,r}\sharp\rho)^{2} \\ &\leq \int_{\mathbb{R}^{d}\times\mathbb{R}^{d}} \|\boldsymbol{x}-\boldsymbol{y}\|^{2} \left[\dot{T}_{k,r}\otimes\dot{T}_{k,r}\right]\sharp\pi_{r}(\mathrm{d}\boldsymbol{x},\mathrm{d}\boldsymbol{y}) \\ &= \int_{\mathbb{R}^{d}\times\mathbb{R}^{d}} \|\dot{T}_{k,r}(\boldsymbol{x})-\dot{T}_{k,r}(\boldsymbol{y})\|^{2}\pi_{r}(\mathrm{d}\boldsymbol{x},\mathrm{d}\boldsymbol{y}) \\ &= \rho(\mathcal{X}_{r})\int_{\mathbb{R}^{d}} \|\dot{T}_{k,r}(\boldsymbol{x})-\dot{T}_{k,r}(\boldsymbol{x})\|^{2}\hat{\mu}_{r}(\mathrm{d}\boldsymbol{x}) \\ &+ (1-\rho(\mathcal{X}_{r}))\int_{\mathbb{R}^{d}\times\mathbb{R}^{d}} \|\dot{T}_{k,r}(\boldsymbol{x})-\dot{T}_{k,r}(\boldsymbol{y})\|^{2}\pi_{r,2}(\mathrm{d}\boldsymbol{x},\mathrm{d}\boldsymbol{y}) \\ &\leq (1-\rho(\mathcal{X}_{r}))\int_{\mathbb{R}^{d}\times\mathbb{R}^{d}} 2\|\dot{T}_{k,r}(\boldsymbol{x})-\dot{T}_{k,r}(\boldsymbol{0})\|^{2}+2\|\dot{T}_{k,r}(\boldsymbol{y})-\dot{T}_{k,r}(\boldsymbol{0})\|^{2}\pi_{r,2}(\mathrm{d}\boldsymbol{x},\mathrm{d}\boldsymbol{y}) \\ &= (1-\rho(\mathcal{X}_{r}))\int_{\mathbb{R}^{d}}2\|\dot{T}_{k,r}(\boldsymbol{x})-\dot{T}_{k,r}(\boldsymbol{0})\|^{2}\hat{\mu}_{r}(\mathrm{d}\boldsymbol{x}) \\ &+ (1-\rho(\mathcal{X}_{r}))\int_{\mathbb{R}^{d}}2\|\dot{T}_{k,r}(\boldsymbol{y})-\dot{T}_{k,r}(\boldsymbol{0})\|^{2}\check{\mu}_{r}(\mathrm{d}\boldsymbol{y}) \\ &\leq \frac{1-\rho(\mathcal{X}_{r})}{\rho(\mathcal{X}_{r})}\int_{\mathbb{R}^{d}}2\|\dot{T}_{k,r}(\boldsymbol{x})-\dot{T}_{k,r}(\boldsymbol{0})\|^{2}\rho(\mathrm{d}\boldsymbol{x})+\int_{\mathbb{R}^{d}}2\|\dot{T}_{k,r}(\boldsymbol{y})-\dot{T}_{k,r}(\boldsymbol{0})\|^{2}\mathbb{1}_{\mathcal{X}_{r}^{c}}(\boldsymbol{y})\rho(\mathrm{d}\boldsymbol{y}). \end{aligned}$$

For k = 1, ..., K, observe that the growth condition of $\dot{T}_{k,r}$ guarantees $\mathbb{E}\left[\left\|\dot{T}_{k,r}(\boldsymbol{x}) - \dot{T}_{k,r}(\boldsymbol{0})\right\|^2\right] \leq u_0(\nu_k) + u_1(\nu_k) \|\boldsymbol{x}\|^2$ for all $\boldsymbol{x} \in \mathbb{R}^d$, where $u_0(\nu_k) \in \mathbb{R}_+$ and $u_1(\nu_k) \in \mathbb{R}_+$ only depend on ν_k . Let $\overline{u}_0 := \max_{1 \leq k \leq K} \{u_0(\nu_k)\}$ and $\overline{u}_1 := \max_{1 \leq k \leq K} \{u_1(\nu_k)\}$. It thus holds that

$$\mathbb{E}\left[\left\|\dot{T}_{k,r}(\boldsymbol{x}) - \dot{T}_{k,r}(\boldsymbol{0})\right\|^{2}\right] \leq \overline{u}_{0} + \overline{u}_{1}\|\boldsymbol{x}\|^{2} \qquad \forall \boldsymbol{x} \in \mathbb{R}^{d}, \ \forall 1 \leq k \leq K.$$
(3.15)

Taking expectations on both sides of (3.14) then applying Fubini's theorem and (3.15) yields

$$\mathbb{E}\left[\mathcal{W}_{2}\left(\dot{T}_{k,r}\sharp\hat{\mu}_{r},\dot{T}_{k,r}\sharp\rho\right)^{2}\right] \\
\leq \frac{1-\rho(\mathcal{X}_{r})}{\rho(\mathcal{X}_{r})}\int_{\mathbb{R}^{d}}2\mathbb{E}\left[\left\|\dot{T}_{k,r}(\boldsymbol{x})-\dot{T}_{k,r}(\boldsymbol{0})\right\|^{2}\right]\rho(\mathrm{d}\boldsymbol{x})+\int_{\mathbb{R}^{d}}2\mathbb{E}\left[\left\|\dot{T}_{k,r}(\boldsymbol{y})-\dot{T}_{k,r}(\boldsymbol{0})\right\|^{2}\right]\mathbb{1}_{\mathcal{X}_{r}^{c}}(\boldsymbol{y})\rho(\mathrm{d}\boldsymbol{y}) \\
\leq \frac{1-\rho(\mathcal{X}_{r})}{\rho(\mathcal{X}_{r})}\int_{\mathbb{R}^{d}}2\left(\overline{u}_{0}+\overline{u}_{1}\|\boldsymbol{x}\|^{2}\right)\rho(\mathrm{d}\boldsymbol{x})+\int_{\mathbb{R}^{d}}2\left(\overline{u}_{0}+\overline{u}_{1}\|\boldsymbol{y}\|^{2}\right)\mathbb{1}_{\mathcal{X}_{r}^{c}}(\boldsymbol{y})\rho(\mathrm{d}\boldsymbol{y}) \qquad \forall 1\leq k\leq K.$$

Same as in the proof of statement (i), since $\bigcup_{r \in \mathbb{N}_0} \mathcal{X}_r = \mathbb{R}^d$ and $\rho \in \mathcal{P}_2(\mathbb{R}^d)$ by assumption, it follows from Lebesgue's dominated convergence theorem that

$$\begin{split} \limsup_{r \to \infty} \mathbb{E} \Big[\mathcal{W}_2 \big(\dot{T}_{k,r} \sharp \hat{\mu}_r, \dot{T}_{k,r} \sharp \rho \big)^2 \Big] &\leq \limsup_{r \to \infty} \frac{1 - \rho(\mathcal{X}_r)}{\rho(\mathcal{X}_r)} \int_{\mathbb{R}^d} 2 \big(\overline{u}_0 + \overline{u}_1 \| \boldsymbol{x} \|^2 \big) \, \rho(\mathrm{d} \boldsymbol{x}) \\ &+ \limsup_{r \to \infty} \int_{\mathbb{R}^d} 2 \big(\overline{u}_0 + \overline{u}_1 \| \boldsymbol{y} \|^2 \big) \, \mathbb{1}_{\mathcal{X}_r^c}(\boldsymbol{y}) \, \rho(\mathrm{d} \boldsymbol{y}) \\ &= 0 \qquad \forall 1 \leq k \leq K. \end{split}$$

Therefore, for any $\epsilon > 0$, there exists $\overline{r}_2(\rho, \nu_1, \dots, \nu_K, \epsilon) \in \mathbb{N}$ such that for any $m_k \ge \underline{m}, n_k \ge \underline{n}, \theta_k \in \Theta$, the inequality $\mathbb{E}\left[\mathcal{W}_2(\dot{T}_{k,r} \sharp \hat{\mu}_r, \dot{T}_{k,r} \sharp \rho)^2\right] \le \epsilon$ holds for all $k = 1, \dots, K$ and all $r \ge \overline{r}_2(\rho, \nu_1, \dots, \nu_K, \epsilon)$. Furthermore, repeating the same derivation in (3.14) with $\dot{T}_{k,r}$ replaced by \dot{T}_r yields

$$\mathcal{W}_{2}\left(\dot{T}_{r}\sharp\hat{\mu}_{r},\dot{T}_{r}\sharp\rho\right)^{2} \leq \frac{1-\rho(\mathcal{X}_{r})}{\rho(\mathcal{X}_{r})} \int_{\mathbb{R}^{d}} 2\left\|\dot{T}_{r}(\boldsymbol{x})-\dot{T}_{r}(\boldsymbol{0})\right\|^{2} \rho(\mathrm{d}\boldsymbol{x}) + \int_{\mathbb{R}^{d}} 2\left\|\dot{T}_{r}(\boldsymbol{y})-\dot{T}_{r}(\boldsymbol{0})\right\|^{2} \mathbb{1}_{\mathcal{X}_{r}^{c}}(\boldsymbol{y}) \rho(\mathrm{d}\boldsymbol{y}).$$

$$(3.16)$$

Observe that, by the convexity of $\mathbb{R}^d \ni \mathbf{z} \mapsto \|\mathbf{z}\|^2 \in \mathbb{R}$, Jensen's inequality, and the growth condition of $(\dot{T}_{k,r})_{k=1:K}$, it holds that

$$\mathbb{E}\Big[\left\|\dot{T}_{r}(\boldsymbol{x})-\dot{T}_{r}(\boldsymbol{0})\right\|^{2}\Big] \leq \sum_{k=1}^{K} w_{k} \mathbb{E}\Big[\left\|\dot{T}_{k,r}(\boldsymbol{x})-\dot{T}_{k,r}(\boldsymbol{0})\right\|^{2}\Big] \leq \overline{u}_{0}+\overline{u}_{1}\|\boldsymbol{x}\|^{2} \qquad \forall \boldsymbol{x} \in \mathbb{R}^{d}.$$
(3.17)

Taking expectations on both sides of (3.16) then applying Fubini's theorem and (3.17) then leads to

$$\mathbb{E}\Big[\mathcal{W}_2\big(\dot{T}_r \sharp \hat{\mu}_r, \dot{T}_r \sharp \rho\big)^2\Big] \leq \frac{1 - \rho(\mathcal{X}_r)}{\rho(\mathcal{X}_r)} \int_{\mathbb{R}^d} 2\big(\overline{u}_0 + \overline{u}_1 \|\boldsymbol{x}\|^2\big) \,\rho(\mathrm{d}\boldsymbol{x}) + \int_{\mathbb{R}^d} 2\big(\overline{u}_0 + \overline{u}_1 \|\boldsymbol{y}\|^2\big) \mathbb{1}_{\mathcal{X}_r^c}(\boldsymbol{y}) \,\rho(\mathrm{d}\boldsymbol{y}).$$

Consequently, it follows from the same argument as above that $\mathbb{E}\left[\mathcal{W}_2(\dot{T}_r \sharp \hat{\mu}_r, \dot{T}_r \sharp \rho)^2\right] \leq \epsilon$ holds whenever $m_k \geq \underline{m}, n_k \geq \underline{n}, \theta_k \in \Theta$, and $r \geq \overline{r}_2(\rho, \nu_1, \dots, \nu_K, \epsilon)$. The proof is now complete.

The results in Proposition 3.10 and Lemma 3.12 suggest the following sufficient conditions for the convergence of Algorithm 2.

Setting 3.13 (Conditions for the convergence of Algorithm 2). Let us fix an arbitrary $\beta > 0$. In addition to Setting 3.6, let the $(\mathcal{F}_t)_{t \in \mathbb{N}_0}$ -adapted stochastic processes $(\widehat{R}_t)_{t \in \mathbb{N}_0}$, $(\widehat{N}_{t,k})_{k=1:K, t \in \mathbb{N}_0}$, and $(\widehat{\Theta}_{t,k})_{k=1:K, t \in \mathbb{N}_0}$ in Algorithm 2 be specified as follows.

(a) For every $t \in \mathbb{N}_0$, let \widehat{R}_t be set as follows:

$$R_{0} := \overline{r}_{2}(\widehat{\rho}_{0}, \nu_{1}, \dots, \nu_{K}, 1),$$

$$\widehat{R}_{t} := \max\left\{\overline{r}_{1}(\widehat{\rho}_{t}, t^{-(1+\beta)}), \overline{r}_{2}(\widehat{\rho}_{t}, \nu_{1}, \dots, \nu_{K}, (t+1)^{-2(1+\beta)})\right\} \qquad \forall t \ge 1$$

where $\overline{r}_1(\cdot, \cdot)$ and $\overline{r}_2(\cdot, \ldots, \cdot)$ are given by Lemma 3.12. Note that \widehat{R}_t is \mathcal{F}_t -measurable for all $t \in \mathbb{N}_0$. (b) For every $t \in \mathbb{N}_0$ and for $k = 1, \ldots, K$, let us specify

$$\begin{split} \widehat{N}_{t,k} &:= \overline{n} \big(\widehat{\mu}_t, \nu_k, (t+1)^{-2(1+\beta)} \big), \\ \widehat{\Theta}_{t,k} &:= \widetilde{\theta} \big(\widehat{\mu}_t, \nu_k, \widehat{N}_{t,k}, \widehat{N}_{t,k}, (t+1)^{-2(1+\beta)} \big) \end{split}$$

where $\overline{n}(\cdot, \cdot, \cdot)$ and $\tilde{\theta}(\cdot, \cdot, \cdot, \cdot, \cdot)$ are given by Assumption 3.4(iii). Note that $\hat{N}_{t,k}$ and $\hat{\Theta}_{t,k}$ are \mathcal{F}_t -measurable for all $t \in \mathbb{N}_0$.

We are now ready to present our main convergence result.

Theorem 3.14 (Convergence of Algorithm 2). Let the inputs of Algorithm 2 satisfy Setting 3.6. Let $(\Omega, \mathcal{F}, \mathbb{P}, (\mathcal{F}_t)_{t \in \mathbb{N}_0})$ be the filtered probability space constructed by Algorithm 2, let the $(\mathcal{F}_t)_{t \in \mathbb{N}_0}$ -adapted stochastic processes $(\widehat{R}_t)_{t \in \mathbb{N}_0}$, $(\widehat{N}_{t,k})_{k=1:K, t \in \mathbb{N}_0}$, and $(\widehat{\Theta}_{t,k})_{k=1:K, t \in \mathbb{N}_0}$ in Algorithm 2 be specified by Setting 3.13, and let $(\widehat{\mu}_t)_{t \in \mathbb{N}_0}$ be the output of Algorithm 2. Then, the following statements hold.

- (i) It holds \mathbb{P} -almost surely that $(\widehat{\mu}_t)_{t \in \mathbb{N}_0}$ is precompact with respect to the \mathcal{W}_2 -metric. Moreover, every accumulation point of $(\widehat{\mu}_t)_{t \in \mathbb{N}_0}$ with respect to the \mathcal{W}_2 -metric is a fixed-point of G.
- (ii) In particular, if G has a unique fixed-point, then $(\widehat{\mu}_t)_{t \in \mathbb{N}_0}$ converges \mathbb{P} -almost surely in \mathcal{W}_2 to the Wasserstein barycenter of ν_1, \ldots, ν_K with weights w_1, \ldots, w_K .

Proof of Theorem 3.14. Let us denote $\overline{T}_{t+1} := \sum_{k=1}^{K} w_k \widehat{T}_{t+1,k}$. Recall that $\widehat{\rho}_{t+1} := \overline{T}_{t+1} \sharp \widehat{\rho}_t$ by Line 10 of Algorithm 2. As implied by Setting 3.13, the properties of $\overline{r}_1(\cdot, \cdot)$, $\overline{r}_2(\cdot, \ldots, \cdot)$ in Lemma 3.12, and the properties of $\overline{n}(\cdot, \cdot, \cdot)$, $\widetilde{\theta}(\cdot, \cdot, \cdot, \cdot, \cdot)$ in Assumption 3.4(iii), the following inequalities hold \mathbb{P} -almost surely:

$$\mathcal{W}_2(\widehat{\mu}_{t+1}, \widehat{\rho}_{t+1})^2 \le (t+1)^{-(1+\beta)} \qquad \forall t \in \mathbb{N}_0, \tag{3.18}$$

$$\mathbb{E}\Big[\mathcal{W}_2\big(\widehat{T}_{t+1,k}\sharp\widehat{\mu}_t,\widehat{T}_{t+1,k}\sharp\widehat{\rho}_t\big)^2\Big|\mathcal{F}_t\Big] \le (t+1)^{-2(1+\beta)} \qquad \forall 1 \le k \le K, \ \forall t \in \mathbb{N}_0,$$
(3.19)

$$\mathbb{E}\Big[\mathcal{W}_2\big(\bar{T}_{t+1}\sharp\hat{\mu}_t,\hat{\rho}_{t+1}\big)^2\Big|\mathcal{F}_t\Big] \le (t+1)^{-2(1+\beta)} \qquad \forall t \in \mathbb{N}_0, \tag{3.20}$$

$$\mathbb{E}\Big[\left\|\widehat{T}_{t+1,k} - T^{\widehat{\mu}_t}_{\nu_k}\right\|_{\mathcal{L}^2(\widehat{\mu}_t)}^2 \Big| \mathcal{F}_t\Big] \le (t+1)^{-2(1+\beta)} \qquad \forall 1 \le k \le K, \ \forall t \in \mathbb{N}_0, \tag{3.21}$$

where $\beta > 0$ is an arbitrary constant. The proof of statement (i) is divided into four steps.

<u>Step 1</u>: showing that $\lim_{t\to\infty} W_2(\widehat{T}_{t+1,k} \sharp \widehat{\rho}_t, \nu_k) = 0 \mathbb{P}$ -almost surely for $k = 1, \ldots, K$. Notice that, for each $t \in \mathbb{N}_0$ and for $k = 1, \ldots, K$, it holds that $[\widehat{T}_{t+1,k}, T_{\nu_k}^{\widehat{\mu}_t}] \sharp \widehat{\mu}_t \in \Pi(\widehat{T}_{t+1,k} \sharp \widehat{\mu}_t, \nu_k)$. Thus, we have

$$\begin{aligned} \mathcal{W}_{2}(\widehat{T}_{t+1,k}\sharp\widehat{\rho}_{t},\nu_{k})^{2} &\leq \left(\mathcal{W}_{2}(\widehat{T}_{t+1,k}\sharp\widehat{\mu}_{t},\widehat{T}_{t+1,k}\sharp\widehat{\rho}_{t}) + \mathcal{W}_{2}(\widehat{T}_{t+1,k}\sharp\widehat{\mu}_{t},\nu_{k})\right)^{2} \\ &\leq 2\mathcal{W}_{2}(\widehat{T}_{t+1,k}\sharp\widehat{\mu}_{t},\widehat{T}_{t+1,k}\sharp\widehat{\rho}_{t})^{2} + 2\mathcal{W}_{2}(\widehat{T}_{t+1,k}\sharp\widehat{\mu}_{t},\nu_{k})^{2} \\ &\leq 2\mathcal{W}_{2}(\widehat{T}_{t+1,k}\sharp\widehat{\mu}_{t},\widehat{T}_{t+1,k}\sharp\widehat{\rho}_{t})^{2} + 2\int_{\mathbb{R}^{d}}\left\|\widehat{T}_{t+1,k}(\boldsymbol{x}) - T_{\nu_{k}}^{\widehat{\mu}_{t}}(\boldsymbol{x})\right\|^{2}\widehat{\mu}_{t}(\boldsymbol{x}) \\ &= 2\mathcal{W}_{2}(\widehat{T}_{t+1,k}\sharp\widehat{\mu}_{t},\widehat{T}_{t+1,k}\sharp\widehat{\rho}_{t})^{2} + 2\left\|\widehat{T}_{t+1,k} - T_{\nu_{k}}^{\widehat{\mu}_{t}}\right\|_{\mathcal{L}^{2}(\widehat{\mu}_{t})}^{2}.\end{aligned}$$

Taking conditional expectations on both sides with respect to \mathcal{F}_t and then applying (3.19) and (3.21) yields

$$\mathbb{E}\Big[\mathcal{W}_2\big(\widehat{T}_{t+1,k}\sharp\widehat{\rho}_t,\nu_k\big)^2\Big|\mathcal{F}_t\Big] \le 4(t+1)^{-2(1+\beta)} \qquad \forall 1\le k\le K, \ \forall t\in\mathbb{N}_0.$$

Subsequently, applying the law of total expectation and Markov's inequality gives

$$\mathbb{P}\Big[\mathcal{W}_2\big(\widehat{T}_{t+1,k}\sharp\widehat{\rho}_t,\nu_k\big)^2 \ge (t+1)^{-(1+\beta)}\Big] \le (t+1)^{1+\beta}\mathbb{E}\Big[\mathcal{W}_2\big(\widehat{T}_{t+1,k}\sharp\widehat{\rho}_t,\nu_k\big)^2\Big] \le 4(t+1)^{-(1+\beta)}$$
$$\forall 1 \le k \le K, \ \forall t \in \mathbb{N}_0.$$

Since $\sum_{t\in\mathbb{N}_0} 4(t+1)^{-(1+\beta)} < \infty$, we conclude by the Borel–Cantelli lemma that, \mathbb{P} -almost surely, $\mathcal{W}_2(\widehat{T}_{t+1,k}\sharp\widehat{\rho}_t,\nu_k)^2 \leq (t+1)^{-(1+\beta)}$ holds for all but finitely many $t \in \mathbb{N}_0$, and it therefore holds that $\lim_{t\to\infty} \mathcal{W}_2(\widehat{T}_{t+1,k}\sharp\widehat{\rho}_t,\nu_k) = 0 \mathbb{P}$ -almost surely for $k = 1,\ldots,K$.

<u>Step 2</u>: showing that $(\widehat{\mu}_t)_{t\in\mathbb{N}_0}$ is precompact with respect to the \mathcal{W}_2 -metric \mathbb{P} -almost surely. By the property that \mathcal{W}_2 metrizes weak convergence in $\mathcal{P}_2(\mathbb{R}^d)$ (see, e.g., [69, Theorem 6.9]) and Prokhorov's theorem, it holds for $k = 1, \ldots, K$ that $(\widehat{T}_{t+1,k} \sharp \widehat{\rho}_t)_{t\in\mathbb{N}_0}$ is \mathbb{P} -almost surely a tight sequence of probability measures. Let $\eta_t := [\widehat{T}_{t+1,1}, \ldots, \widehat{T}_{t+1,K}] \sharp \widehat{\rho}_t \in \mathcal{P}(\underbrace{\mathbb{R}^d \times \cdots \times \mathbb{R}^d}_{K \text{ copies}})$ for $t \in \mathbb{N}_0$. It hence holds \mathbb{P} -almost surely that each

marginal of the sequence $(\eta_t)_{t\in\mathbb{N}_0}$ (on each copy of \mathbb{R}^d) belongs to a tight set of probability measures on \mathbb{R}^d , and it then follows from a multi-marginal generalization of [69, Lemma 4.4] that $(\eta_t)_{t\in\mathbb{N}_0}$ is a tight set of probability measures on $(\mathbb{R}^d)^K$. Consequently, Prokhorov's theorem implies that every subsequence of $(\eta_t)_{t\in\mathbb{N}_0}$ admits a further subsequence which is weakly convergent. Let $(\eta_{t_i})_{i\in\mathbb{N}_0}$ be a weakly convergent subsequence of $(\eta_t)_{t\in\mathbb{N}_0}$ with weak limit $\eta_{t_{\infty}} \in \mathcal{P}((\mathbb{R}^d)^K)$. It subsequently follows from Step 1 that $\eta_{t_{\infty}} \in \Pi(\nu_1, \ldots, \nu_K)$, and hence

$$\lim_{i \to \infty} \int_{(\mathbb{R}^d)^K} \|\boldsymbol{x}\|^2 \eta_{t_i}(\mathrm{d}\boldsymbol{x}) = \lim_{i \to \infty} \int_{(\mathbb{R}^d)^K} \sum_{k=1}^K \|\boldsymbol{x}_k\|^2 \eta_{t_i}(\mathrm{d}\boldsymbol{x}_1, \dots, \mathrm{d}\boldsymbol{x}_K)$$
$$= \lim_{i \to \infty} \sum_{k=1}^K \int_{\mathbb{R}^d} \|\boldsymbol{x}_k\|^2 \widehat{T}_{t_i+1,k} \sharp \widehat{\rho}_{t_i}(\mathrm{d}\boldsymbol{x}_k)$$
$$= \sum_{k=1}^K \int_{\mathbb{R}^d} \|\boldsymbol{x}_k\|^2 \nu_k(\mathrm{d}\boldsymbol{x}_k) = \int_{(\mathbb{R}^d)^K} \|\boldsymbol{x}\|^2 \eta_{t_\infty}(\mathrm{d}\boldsymbol{x}) \qquad \mathbb{P}\text{-a.s.}$$

Next, let A denote the mapping $(\mathbb{R}^d)^K \ni (\boldsymbol{x}_1, \dots, \boldsymbol{x}_K) \mapsto \sum_{k=1}^K w_k \boldsymbol{x}_k \in \mathbb{R}^d$. Hence, we have $\hat{\rho}_{t+1} = \overline{T}_{t+1} \sharp \hat{\rho}_t = A \sharp \eta_t$ for all $t \in \mathbb{N}_0$. It follows from the convexity of $\mathbb{R}^d \ni \boldsymbol{z} \mapsto \|\boldsymbol{z}\|^2 \in \mathbb{R}$ and Jensen's inequality that

$$\|A(\boldsymbol{x})\|^{2} \leq \sum_{k=1}^{K} w_{k} \|\boldsymbol{x}_{k}\|^{2} \leq \|\boldsymbol{x}\|^{2} \qquad \forall \boldsymbol{x} = (\boldsymbol{x}_{1}, \dots, \boldsymbol{x}_{K}) \in (\mathbb{R}^{d})^{K}.$$
(3.23)

Combining (3.23), (3.22), and the equivalence between (iii) and (iv) in [68, Theorem 7.12] yields

$$\lim_{i \to \infty} \int_{\mathbb{R}^d} \|\boldsymbol{y}\|^2 \, \widehat{\rho}_{t_i+1}(\mathrm{d}\boldsymbol{y}) = \lim_{i \to \infty} \int_{(\mathbb{R}^d)^K} \|A(\boldsymbol{x})\|^2 \, \eta_{t_i}(\mathrm{d}\boldsymbol{x})$$

$$= \int_{(\mathbb{R}^d)^K} \|A(\boldsymbol{x})\|^2 \, \eta_{t_\infty}(\mathrm{d}\boldsymbol{x}) = \int_{\mathbb{R}^d} \|\boldsymbol{y}\|^2 \, A \sharp \eta_{t_\infty}(\mathrm{d}\boldsymbol{y}) \qquad \mathbb{P}\text{-a.s.}$$
(3.24)

Moreover, since $(\eta_{t_i})_{i \in \mathbb{N}_0}$ converges weakly to η_{t_∞} \mathbb{P} -almost surely and A is continuous, it holds \mathbb{P} -almost surely that $\lim_{i\to\infty} \int_{\mathbb{R}^d} \psi \, d\hat{\rho}_{t_i+1} = \lim_{i\to\infty} \int_{(\mathbb{R}^d)^K} \psi \circ A \, d\eta_{t_i} = \int_{(\mathbb{R}^d)^K} \psi \circ A \, d\eta_{t_\infty} = \int_{\mathbb{R}^d} \psi \, dA \sharp \eta_{t_\infty}$ for any continuous and bounded function $\psi : \mathbb{R}^d \to \mathbb{R}$, which shows that $(\hat{\rho}_{t_i+1})_{i\in\mathbb{N}_0}$ converges weakly to $A\sharp\eta_{t_\infty}$ \mathbb{P} -almost surely. Now, (3.24) and the equivalence between (i) and (iii) in [68, Theorem 7.12] show that $\lim_{i\to\infty} \mathcal{W}_2(\hat{\rho}_{t_i+1}, A\sharp\eta_{t_\infty}) = 0$ \mathbb{P} -almost surely. Furthermore, (3.18) implies that $\lim_{t\to\infty} \mathcal{W}_2(\hat{\mu}_t, \hat{\rho}_t) = 0$ \mathbb{P} -almost surely. The above analyses have established that, \mathbb{P} -almost surely, every subsequence of $(\hat{\mu}_t)_{t\in\mathbb{N}_0}$ admits a further subsequence which converges with respect to the \mathcal{W}_2 -metric, and thus $(\hat{\mu}_t)_{t\in\mathbb{N}_0}$ is precompact with respect to the \mathcal{W}_2 -metric \mathbb{P} -almost surely.

<u>Step 3</u>: constructing an \mathcal{F} -measurable set $\Omega \subseteq \Omega$ with $\mathbb{P}[\Omega] = 1$ in which the convergence is analyzed. Similar to the argument used in Step 1, applications of the law of total expectation together with Markov's inequality to (3.20) and (3.21) lead to

$$\mathbb{P}\Big[\mathcal{W}_{2}\big(\bar{T}_{t+1} \sharp \hat{\mu}_{t}, \hat{\rho}_{t+1}\big)^{2} \ge (t+1)^{-(1+\beta)}\Big] \le (t+1)^{-(1+\beta)} \qquad \forall t \in \mathbb{N}_{0},$$
$$\mathbb{P}\Big[\big\|\widehat{T}_{t+1,k} - T_{\nu_{k}}^{\widehat{\mu}_{t}}\big\|_{\mathcal{L}^{2}(\widehat{\mu}_{t})}^{2} \ge (t+1)^{-(1+\beta)}\Big] \le (t+1)^{-(1+\beta)} \qquad \forall 1 \le k \le K, \ \forall t \in \mathbb{N}_{0},$$

Since $\sum_{t\in\mathbb{N}_0} (t+1)^{-(1+\beta)} < \infty$, we use the Borel–Cantelli lemma again to show that, \mathbb{P} -almost surely, $\mathcal{W}_2(\bar{T}_{t+1}\sharp\hat{\mu}_t,\hat{\rho}_{t+1})^2 \leq (t+1)^{-(1+\beta)}$ and $\|\hat{T}_{t+1,k} - T_{\nu_k}^{\hat{\mu}_t}\|_{\mathcal{L}^2(\hat{\mu}_t)}^2 \leq (t+1)^{-(1+\beta)} \forall 1 \leq k \leq K$ hold for all but finitely many $t \in \mathbb{N}_0$. In the following, for every $\omega \in \Omega$, let us use the notations $\hat{\rho}_t^{(\omega)}, \hat{\mu}_t^{(\omega)}, \hat{T}_{t+1,k}^{(\omega)}, \bar{T}_{t+1}^{(\omega)}$ to explicitly express the dependence of the random variables $\hat{\rho}_t, \hat{\mu}_t, \hat{T}_{t+1,k}, \bar{T}_{t+1}$ on ω . The above analyses have shown the existence of an \mathcal{F} -measurable set $\tilde{\Omega} \subseteq \Omega$ with $\mathbb{P}[\tilde{\Omega}] = 1$, which satisfies:

$$\forall \omega \in \widetilde{\Omega}, \ \exists \overline{t}^{(\omega)} \in \mathbb{N}_{0}, \begin{cases} \left(\widehat{\mu}_{t}^{(\omega)}\right)_{t \in \mathbb{N}_{0}} \text{ is precompact with respect to the } \mathcal{W}_{2}\text{-metric,} \\ \mathcal{W}_{2}\left(\widehat{\mu}_{t+1}^{(\omega)}, \widehat{\rho}_{t+1}^{(\omega)}\right)^{2} \leq (t+1)^{-(1+\beta)} & \forall t \geq \overline{t}^{(\omega)}, \\ \mathcal{W}_{2}\left(\overline{T}_{t+1}^{(\omega)} \# \widehat{\mu}_{t}^{(\omega)}, \widehat{\rho}_{t+1}^{(\omega)}\right)^{2} \leq (t+1)^{-(1+\beta)} & \forall t \geq \overline{t}^{(\omega)}, \\ \left\|\widehat{T}_{t+1,k}^{(\omega)} - T_{\nu_{k}}^{\widehat{\mu}_{t}^{(\omega)}}\right\|_{\mathcal{L}^{2}\left(\widehat{\mu}_{t}^{(\omega)}\right)}^{2} \leq (t+1)^{-(1+\beta)} & \forall 1 \leq k \leq K, \ \forall t \geq \overline{t}^{(\omega)}. \end{cases}$$
(3.25)

<u>Step 4</u>: showing that for every $\omega \in \widetilde{\Omega}$, every \mathcal{W}_2 -accumulation point of $(\widehat{\mu}_t^{(\omega)})_{t\in\mathbb{N}_0}$ is a fixed-point of G. Let us fix an arbitrary $\omega \in \widetilde{\Omega}$ and let the subsequence $(t_i)_{i\in\mathbb{N}_0}$ be such that $\lim_{i\to\infty} \mathcal{W}_2(\widehat{\mu}_{t_i}^{(\omega)}, \widehat{\mu}_{\infty}^{(\omega)}) = 0$ for $\widehat{\mu}_{\infty}^{(\omega)} \in \mathcal{P}_2(\mathbb{R}^d)$. The continuity of $V(\cdot)$ on $\mathcal{P}_2(\mathbb{R}^d)$ then implies that $\lim_{i\to\infty} V(\widehat{\mu}_{t_i}^{(\omega)}) = V(\widehat{\mu}_{\infty}^{(\omega)})$. Removing finitely many initial terms from $(t_i)_{i\in\mathbb{N}_0}$ if necessary, we assume without loss of generality that $t_0 \geq \overline{t}^{(\omega)}$. For each $i \in \mathbb{N}_0$, summing (3.3) over $s = t_i, t_i + 1, \ldots, t_{i+1} - 1$, using the inequality $\mathcal{W}_2(\overline{T}_{s+1}^{(\omega)} \sharp \widehat{\mu}_s^{(\omega)}, \widehat{\mu}_{s+1}^{(\omega)})^2 \leq 2\mathcal{W}_2(\overline{T}_{s+1}^{(\omega)} \sharp \widehat{\mu}_s^{(\omega)}, \widehat{\rho}_{s+1}^{(\omega)})^2 + 2\mathcal{W}_2(\widehat{\mu}_{s+1}^{(\omega)}, \widehat{\rho}_{s+1}^{(\omega)})^2$, and using the properties in (3.25) lead to

$$\begin{split} V(\widehat{\mu}_{t_{i+1}}^{(\omega)}) - V(\widehat{\mu}_{t_{i}}^{(\omega)}) &= \sum_{s=t_{i}}^{t_{i+1}-1} V(\widehat{\mu}_{s+1}^{(\omega)}) - V(\widehat{\mu}_{s}^{(\omega)}) \\ &\leq -\left(\sum_{s=t_{i}}^{t_{i+1}-1} \mathcal{W}_{2}\Big(\widehat{\mu}_{s}^{(\omega)}, G(\widehat{\mu}_{s}^{(\omega)})\Big)^{2}\right) + \left(\sum_{s=t_{i}}^{t_{i+1}-1} 2\sum_{k=1}^{K} w_{k} \left\|\widehat{T}_{s+1,k}^{(\omega)} - T_{\nu_{k}}^{\widehat{\mu}_{s}^{(\omega)}}\right\|_{\mathcal{L}^{2}(\widehat{\mu}_{s}^{(\omega)})}\right) \\ &+ \left(\sum_{s=t_{i}}^{t_{i+1}-1} 4\mathcal{W}_{2}(\overline{T}_{s+1}^{(\omega)} + \widehat{\mu}_{s}^{(\omega)}, \widehat{\rho}_{s+1}^{(\omega)})^{2}\right) + \left(\sum_{s=t_{i}}^{t_{i+1}-1} 4\mathcal{W}_{2}(\widehat{\mu}_{s+1}^{(\omega)}, \widehat{\rho}_{s+1}^{(\omega)})^{2}\right) \\ &\leq -\left(\sum_{s=t_{i}}^{t_{i+1}-1} \mathcal{W}_{2}(\widehat{\mu}_{s}^{(\omega)}, G(\widehat{\mu}_{s}^{(\omega)})\Big)^{2}\right) + \left(\sum_{s=t_{i}}^{t_{i+1}-1} 10(s+1)^{-(1+\beta)}\right) \\ &\leq -\mathcal{W}_{2}(\widehat{\mu}_{t_{i}}^{(\omega)}, G(\widehat{\mu}_{t_{i}}^{(\omega)})\Big)^{2} + \left(\sum_{s=t_{i}}^{\infty} 10(s+1)^{-(1+\beta)}\right) \qquad \forall i \in \mathbb{N}_{0}. \end{split}$$

Rearranging the terms above leads to

$$\mathcal{W}_2\left(\widehat{\mu}_{t_i}^{(\omega)}, G(\widehat{\mu}_{t_i}^{(\omega)})\right)^2 \le \left| V(\widehat{\mu}_{t_{i+1}}^{(\omega)}) - V(\widehat{\mu}_{t_i}^{(\omega)}) \right| + \left(\sum_{s=t_i}^{\infty} 10(s+1)^{1+\beta}\right) \qquad \forall i \in \mathbb{N}_0$$

Since $\sum_{s=0}^\infty (s+1)^{-(1+\beta)}$ is a convergent series, we get

$$\limsup_{i \to \infty} \mathcal{W}_2\left(\widehat{\mu}_{t_i}^{(\omega)}, G(\widehat{\mu}_{t_i}^{(\omega)})\right)^2 \le \limsup_{i \to \infty} \left| V(\widehat{\mu}_{t_{i+1}}^{(\omega)}) - V(\widehat{\mu}_{t_i}^{(\omega)}) \right| + \limsup_{i \to \infty} \left(\sum_{s=t_i}^{\infty} 10(s+1)^{1+\beta} \right) = 0.$$

This shows that $\lim_{i\to\infty} W_2(\widehat{\mu}_{t_i}^{(\omega)}, G(\widehat{\mu}_{t_i}^{(\omega)})) = 0$. Moreover, for any $\mu \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$, the analysis in [3, Remark 3.2] demonstrates that the density function $f_{G(\mu)}$ of $G(\mu) \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$ satisfies

$$\sup_{\boldsymbol{x}\in\mathbb{R}^d}\left\{f_{G(\mu)}(\boldsymbol{x})\right\}\leq w_1^{-d}\sup_{\boldsymbol{x}\in\mathrm{supp}(\nu_1)}\left\{f_{\nu_1}(\boldsymbol{x})\right\}<\infty,$$

where f_{ν_1} denotes the density function of $\nu_1 \in \mathcal{M}^q(\mathbb{R}^d)$. Consequently, it holds for every open set $E \subseteq \mathbb{R}^d$ that

$$\widehat{\mu}_{\infty}^{(\omega)}(E) \leq \liminf_{i \to \infty} G(\widehat{\mu}_{t_i}^{(\omega)})(E) \leq w_1^{-d} \sup_{\boldsymbol{x} \in \operatorname{supp}(\nu_1)} \{f_{\nu_1}(\boldsymbol{x})\} \mathscr{L}(E),$$

where \mathscr{L} denotes the Lebesgue measure on \mathbb{R}^d . It thus follows that $\widehat{\mu}_{\infty}^{(\omega)} \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$. Now, the continuity of the mapping $\mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d) \ni \mu \mapsto \mathcal{W}_2(\mu, G(\mu))^2 \in \mathbb{R}_+$ [3, Theorem 3.1] implies that $\mathcal{W}_2(\widehat{\mu}_{\infty}^{(\omega)}, G(\widehat{\mu}_{\infty}^{(\omega)}))^2 = \lim_{i\to\infty} \mathcal{W}_2(\widehat{\mu}_{t_i}^{(\omega)}, G(\widehat{\mu}_{t_i}^{(\omega)}))^2 = 0$, which shows that $\widehat{\mu}_{\infty}^{(\omega)}$ is a fixed-point of G. Since $\mathbb{P}[\widetilde{\Omega}] = 1$, it holds \mathbb{P} -almost surely that every \mathcal{W}_2 -accumulation point of $(\widehat{\mu}_t)_{t\in\mathbb{N}_0}$ is a fixed-point of G. We have thus completed the proof of statement (i).

Finally, if G has a unique fixed-point $\bar{\mu} \in \mathcal{P}_{2,ac}(\mathbb{R}^d)$, then statement (i) implies that, \mathbb{P} -almost surely, every \mathcal{W}_2 -accumulation point of $(\hat{\mu}_t)_{t \in \mathbb{N}_0}$ is equal to $\bar{\mu}$. Therefore, $(\hat{\mu}_t)_{t \in \mathbb{N}_0}$ converges \mathbb{P} -almost surely in \mathcal{W}_2 to $\bar{\mu}$, which is the unique Wasserstein barycenter of ν_1, \ldots, ν_K by Theorem 1.2(i). The proof is now complete. \Box

Remark 3.15. We would like to remark that the operator G in (1.3) does not always have a unique fixedpoint for general input probability measures $\nu_1, \ldots, \nu_K \in \mathcal{P}_{2,ac}(\mathbb{R}^d)$; see, e.g., Example 3.1 of [3] for a concrete counterexample. It is known that G has a unique fixed-point when ν_1, \ldots, ν_K belong to the same parametric family of elliptical distributions [3, Section 4], e.g., Gaussian distributions. However, to the best of our knowledge, sufficient conditions to guarantee the uniqueness of the fixed-point of G for non-parametric ν_1, \ldots, ν_K is still an open problem.

Remark 3.16. Same as the deterministic fixed-point iterative scheme of Álvarez-Esteban et al. [3], our stochastic extension in Algorithm 2 does not provide a rate of convergence. It will be shown in our numerical experiments (see Section 6 for details), however, that our algorithm tends to converge right after the first few iterations empirically, which coincides with an analogous empirical observation of von Lindheim [70] in a setting involving only discrete measures.

4. MODIFIED ENTROPIC OT MAP ESTIMATOR

As stated in Setting 3.13, the convergence of Algorithm 2 depends crucially on the OT map estimator $\widehat{T}_{\nu,n}^{\mu,m}[\theta]$, specifically on its shape, growth, and consistency properties required by Assumption 3.4. In this section, we consider two admissible compactly supported probability measures $\mu, \nu \in \mathcal{M}^q(\mathbb{R}^d)$ for $q \ge 1$ and introduce a concrete example of OT map estimator that satisfies Assumption 3.4, which is a modified version of the entropic OT map estimator of Pooladian and Niles-Weed [52]. This estimator is constructed via solving an entropic optimal transport problem [20] followed by the operation of barycentric projection [4, Definition 5.4.2].

For the sake of notational simplicity, we will omit μ, ν, m, n in the notations and denote the entropic OT map estimator by $\widehat{T}_{entr}[\theta]$. Nonetheless, m and n will always be understood as the numbers of samples from μ and ν , respectively. Our modification, compared to [52], lies in the addition of a term that guarantees the strong

convexity condition in Assumption 3.4(i). The following proposition presents the definition and properties of $\hat{T}_{entr}[\theta]$.

Proposition 4.1 (Modified Entropic OT map estimator). Let $q \ge 1$, $\mu, \nu \in \mathcal{M}^q(\mathbb{R}^d)$ (recall Definition 3.1), $\underline{m} := 1$, $\underline{n} := d + 1$, $\Theta := (0, \infty)$, let T^{μ}_{ν} be the OT map from μ to ν , and let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space. Moreover, let $r_0(\mu) := \inf \{r \in \mathbb{R}_+ : \operatorname{supp}(\mu) \subseteq \overline{B}(\mathbf{0}, r)\}$, $r_0(\nu) := \inf \{r \in \mathbb{R}_+ : \operatorname{supp}(\nu) \subseteq \overline{B}(\mathbf{0}, r)\}$, and let $\gamma : \mathbb{R}_+ \to \mathbb{R}_+$ be defined as follows:

$$\gamma(z) := \begin{cases} 0 & \text{if } z \le \frac{r_0(\mu)^2}{2}, \\ \exp\left(-\frac{2}{2z - r_0(\mu)^2}\right) & \text{if } z > \frac{r_0(\mu)^2}{2} \end{cases} \quad \forall z \in \mathbb{R}_+.$$

Given $m \geq \underline{m}$ independent random samples $\mathbf{X}_1, \ldots, \mathbf{X}_m : \Omega \to \mathbb{R}^d$ from μ and $n \geq \underline{n}$ independent random samples $\mathbf{Y}_1, \ldots, \mathbf{Y}_n : \Omega \to \mathbb{R}^d$ from ν , we construct $\widehat{T}_{entr}[\theta] : \mathbb{R}^d \to \mathbb{R}^d$ through two steps.

(1) Sinkhorn step: for every $\theta > 0$, let $(\widehat{f}_i^{\theta})_{i=1:m}$, $(\widehat{g}_j^{\theta})_{j=1:n}$ be the optimal solution of the dual entropic optimal transport problem:

$$\begin{array}{ll} \underset{(\widehat{f}_i),(\widehat{g}_j)}{\text{maximize}} & \left(\frac{1}{m}\sum_{i=1}^m\widehat{f}_i\right) + \left(\frac{1}{n}\sum_{j=1}^n\widehat{g}_j\right) - \left(\frac{\theta}{mn}\sum_{i=1}^m\sum_{j=1}^n\exp\left(\frac{1}{\theta}\left(\widehat{f}_i + \widehat{g}_j - \frac{1}{2}\|\boldsymbol{X}_i - \boldsymbol{Y}_j\|^2\right)\right)\right) \\ \text{subject to} & \widehat{f}_i \in \mathbb{R} \quad \forall 1 \le i \le m, \qquad \widehat{g}_j \in \mathbb{R} \quad \forall 1 \le j \le n. \end{array}$$

$$(4.1)$$

(2) Barycentric projection step: let $\widehat{\varphi}_{entr}[\theta] : \mathbb{R}^d \to \mathbb{R}$ and $\widehat{T}_{entr}[\theta] : \mathbb{R}^d \to \mathbb{R}^d$ be defined as follows:

$$\widehat{\varphi}_{\text{entr}}[\theta](\boldsymbol{x}) := \theta \log \left(\sum_{j=1}^{n} \exp\left(\frac{1}{\theta} \left(\widehat{g}_{j}^{\theta} + \langle \boldsymbol{Y}_{j}, \boldsymbol{x} \rangle - \frac{1}{2} \|\boldsymbol{Y}_{j}\|^{2} \right) \right) \right) + \int_{0}^{\|\boldsymbol{x}\|^{2}/2} \gamma(z) \, \mathrm{d}z \quad \forall \boldsymbol{x} \in \mathbb{R}^{d}, \quad (4.2)$$

$$\widehat{T}_{entr}[\theta](\boldsymbol{x}) := \frac{\sum_{j=1}^{n} \exp\left(\frac{1}{\theta} \left(\widehat{g}_{j}^{\theta} + \langle \boldsymbol{Y}_{j}, \boldsymbol{x} \rangle - \frac{1}{2} \|\boldsymbol{Y}_{j}\|^{2}\right)\right) \boldsymbol{Y}_{j}}{\sum_{j=1}^{n} \exp\left(\frac{1}{\theta} \left(\widehat{g}_{j}^{\theta} + \langle \boldsymbol{Y}_{j}, \boldsymbol{x} \rangle - \frac{1}{2} \|\boldsymbol{Y}_{j}\|^{2}\right)\right)} + \gamma\left(\frac{1}{2} \|\boldsymbol{x}\|^{2}\right) \boldsymbol{x} \qquad \forall \boldsymbol{x} \in \mathbb{R}^{d}.$$
(4.3)

Then, $\widehat{T}_{entr}[\theta] \in C_{lin}(\mathbb{R}^d, \mathbb{R}^d)$ has a Borel dependence on $(X_1, \ldots, X_m, Y_1, \ldots, Y_n, \theta)$, and the following statements hold.

- (i) Shape: for all $m \geq \underline{m}$, $n \geq \underline{n}$, and for all $\theta > 0$, it holds that $\widehat{T}_{entr}[\theta](\boldsymbol{x}) = \nabla \widehat{\varphi}_{entr}[\theta](\boldsymbol{x})$ for all $\boldsymbol{x} \in \mathbb{R}^d$. Moreover, let $\underline{\lambda}(\mu, \nu, m, n, \boldsymbol{X}_1, \dots, \boldsymbol{X}_m, \boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n, \theta) := \min\left\{\frac{1}{\theta}\exp\left(-\frac{1}{\theta}(6r_0(\mu) + 4r_0(\nu))r_0(\nu)\right)e_{\min}(\widehat{Cov}[\boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n]), \exp\left(-\frac{2}{3r_0(\mu)^2}\right)\right\}$ (which is abbreviated to $\underline{\lambda}$ in the following), where $\widehat{Cov}[\boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n] := \left(\frac{1}{n}\sum_{j=1}^n \boldsymbol{Y}_j\boldsymbol{Y}_j^{\mathsf{T}}\right) - \left(\frac{1}{n}\sum_{j=1}^n \boldsymbol{Y}_j\right)\left(\frac{1}{n}\sum_{j=1}^n \boldsymbol{Y}_j\right)^{\mathsf{T}} \in \mathbb{R}^{d\times d}$ denotes the (biased) sample covariance matrix of $\boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n$. Then, it holds \mathbb{P} -almost surely that $\underline{\lambda} > 0$ and $\widehat{\varphi}_{entr}[\theta] \in \mathfrak{C}_{\underline{\lambda},\infty}^{\infty}(\mathbb{R}^d)$. In particular, $\widehat{T}_{entr}[\theta]$ satisfies Assumption 3.4(i) with respect to $\alpha(\mu, \nu, m, n, \boldsymbol{X}_1, \dots, \boldsymbol{X}_m, \boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n, \theta) \leftarrow 1$ and $\underline{\lambda}(\mu, \nu, m, n, \boldsymbol{X}_1, \dots, \boldsymbol{X}_m, \boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n, \theta) \leftarrow$ $\min\left\{\frac{1}{\theta}\exp\left(-\frac{1}{\theta}(6r_0(\mu) + 4r_0(\nu))r_0(\nu)\right)e_{\min}(\widehat{Cov}[\boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n]), \exp\left(-\frac{2}{3r_0(\mu)^2}\right)\right\}.$
- (ii) Growth: for all $m \geq \underline{m}$, $n \geq \underline{n}$, and all $\theta > 0$, it holds \mathbb{P} -almost surely that $\|\widehat{T}_{entr}[\theta](\boldsymbol{x}) \widehat{T}_{entr}[\theta](\boldsymbol{0})\|^2 \leq 8r_0(\nu)^2 + 2\|\boldsymbol{x}\|^2$ for all $\boldsymbol{x} \in \mathbb{R}^d$. In particular, $\widehat{T}_{entr}[\theta]$ satisfies Assumption 3.4(ii) with respect to $u_0(\nu) \leftarrow 8r_0(\nu)^2$, $u_1(\nu) \leftarrow 2$.
- (iii) Consistency: there exist $\overline{\theta}(\mu,\nu) > 0$, $3 \le \overline{\alpha}(\mu,\nu) \le 4$, and $C_{\text{entr}}(\mu,\nu) > 0$ such that

$$\begin{split} \mathbb{E}\Big[\left\| \widehat{T}_{\text{entr}}[\theta] - T^{\mu}_{\nu} \right\|_{\mathcal{L}^{2}(\mu)}^{2} \Big] &\leq C_{\text{entr}}(\mu,\nu) \Big[\theta^{-\frac{d}{2}} \big(\log(n)n^{-\frac{1}{2}} + \log(m)m^{-\frac{1}{2}} \big) + \theta^{\frac{\overline{\alpha}(\mu,\nu)}{2}} \\ &\forall 0 < \theta \leq \overline{\theta}(\mu,\nu), \; \forall m \geq \underline{m}, \; \forall n \geq \underline{n}. \end{split}$$

In particular, $\hat{T}_{entr}[\theta]$ satisfies Assumption 3.4(iii) with respect to

$$\overline{n}(\mu,\nu,\epsilon) \leftarrow \min\left\{n \in \mathbb{N}: \begin{array}{l} n \geq \max\left\{d+1,7,\overline{\theta}(\mu,\nu)^{-(\overline{\alpha}(\mu,\nu)+d)}\right\},\\ m^{-\frac{\overline{\alpha}(\mu,\nu)}{2(\overline{\alpha}(\mu,\nu)+d)}}\left(\log(m)+1\right) \leq \frac{\epsilon}{C_{\mathrm{entr}}(\mu,\nu)} \ \forall m \geq n \end{array}\right\},\\ \widetilde{\theta}(\mu,\nu,m,n,\epsilon) \leftarrow \min\{m,n\}^{-\frac{1}{\overline{\alpha}(\mu,\nu)+d}}.$$

Before we prove Proposition 4.1, let us first establish the following lemma.

Lemma 4.2. Let $\theta > 0$, $n \in \mathbb{N}$, let $g_j \in \mathbb{R}$, $y_j \in \mathbb{R}^d$ for j = 1, ..., n, and let $\eta_1, ..., \eta_n : \mathbb{R}^d \to [0, 1]$, $\varphi : \mathbb{R}^d \to \mathbb{R}$, $T : \mathbb{R}^d \to \mathbb{R}^d$ be defined as follows:

$$egin{aligned} &\eta_j(m{x}) \coloneqq rac{\exp\left(rac{1}{ heta}ig(g_j+\langlem{y}_j,m{x}
angle-rac{1}{2}\|m{y}_j\|^2ig)
ight)}{\sum_{j'=1}^n \exp\left(rac{1}{ heta}ig(g_{j'}+\langlem{y}_{j'},m{x}
angle-rac{1}{2}\|m{y}_{j'}\|^2ig)
ight)} &oralligxid x\in\mathbb{R}^d, \ orall 1\leq j\leq n, \ &arphi(m{x})\coloneqq heta\log\left(\sum_{j=1}^n \exp\left(rac{1}{ heta}ig(g_j+\langlem{y}_j,m{x}
angle-rac{1}{2}\|m{y}_j\|^2ig)
ight)
ight) &orall igxid x\in\mathbb{R}^d, \ &arphi(x)\leq j\leq n, \ &arphi(m{x})\coloneqq heta\log\left(\sum_{j=1}^n \exp\left(rac{1}{ heta}ig(g_j+\langlem{y}_j,m{x}
angle-rac{1}{2}\|m{y}_j\|^2ig)igright)
ight) &orall igxid x\in\mathbb{R}^d, \ &arphi(m{x}\in\mathbb{R}^d, \ &arphi(m{$$

Then, φ is convex, $\varphi \in \mathcal{C}^{\infty}(\mathbb{R}^d)$, and $\nabla \varphi = T$. Moreover, it holds that

$$\frac{n}{\theta} \min_{1 \le j \le n} \left\{ \eta_j(\boldsymbol{x}) \right\} e_{\min} \Big(\widehat{\operatorname{Cov}}[\boldsymbol{y}_1, \dots, \boldsymbol{y}_n] \Big) \mathbf{I}_d \preceq \nabla^2 \varphi(\boldsymbol{x}) \preceq \frac{1}{\theta} \max_{1 \le j \le n} \left\{ \|\boldsymbol{y}_j\|^2 \right\} \mathbf{I}_d \qquad \forall \boldsymbol{x} \in \mathbb{R}^d$$

where $\widehat{\operatorname{Cov}}[\boldsymbol{y}_1, \dots, \boldsymbol{y}_n] := \left(\frac{1}{n} \sum_{j=1}^n \boldsymbol{y}_j \boldsymbol{y}_j^{\mathsf{T}}\right) - \left(\frac{1}{n} \sum_{j=1}^n \boldsymbol{y}_j\right) \left(\frac{1}{n} \sum_{j=1}^n \boldsymbol{y}_j\right)^{\mathsf{T}} \in \mathbb{R}^{d \times d}$ denotes the (biased) sample covariance matrix of $\boldsymbol{y}_1, \dots, \boldsymbol{y}_n$.

Proof of Lemma 4.2. To begin, let $\Delta_n := \left\{ u = (u_1, \dots, u_n)^{\mathsf{T}} \in \mathbb{R}^n : \sum_{j=1}^n u_j = 1, u_j \ge 0 \ \forall 1 \le j \le n \right\}$. Observe that φ is the composition of a log-sum-exp function and an affine function and is thus convex. Moreover, it can be directly verified from the definitions that $\varphi \in \mathcal{C}^{\infty}(\mathbb{R})$ and $\nabla \varphi = T$. Next, let **Y** denote the matrix formed with y_1, \dots, y_n as columns, that is,

$$\mathbf{Y} := \begin{pmatrix} | & | & | \\ \boldsymbol{y}_1 & \boldsymbol{y}_2 & \cdots & \boldsymbol{y}_n \\ | & | & | \end{pmatrix} \in \mathbb{R}^{d \times n}$$

Notice that

$$abla \eta_j(oldsymbol{x}) = rac{1}{ heta} \eta_j(oldsymbol{x}) oldsymbol{y}_j - rac{1}{ heta} \eta_j(oldsymbol{x}) igg(oldsymbol{x}) \sum_{l=1}^n \eta_l(oldsymbol{x}) oldsymbol{y}_l igg) \qquad orall oldsymbol{x} \in \mathbb{R}^d, \ orall 1 \leq j \leq n.$$

Letting $\boldsymbol{\eta}(\boldsymbol{x}) := \left(\eta_1(\boldsymbol{x}), \dots, \eta_n(\boldsymbol{x})\right)^{\mathsf{T}} \in \Delta_n$, we hence get

$$\nabla^{2}\varphi(\boldsymbol{x}) = \sum_{j=1}^{n} \nabla \eta_{j}(\boldsymbol{x})\boldsymbol{y}_{j}^{\mathsf{T}} = \frac{1}{\theta} \left(\sum_{j=1}^{n} \eta_{j}(\boldsymbol{x})\boldsymbol{y}_{j}\boldsymbol{y}_{j}^{\mathsf{T}} \right) - \frac{1}{\theta} \left(\sum_{j=1}^{n} \eta_{j}(\boldsymbol{x})\boldsymbol{y}_{j} \right) \left(\sum_{j=1}^{n} \eta_{j}(\boldsymbol{x})\boldsymbol{y}_{j} \right)^{\mathsf{T}}$$

$$= \frac{1}{\theta} \mathbf{Y} \left(\operatorname{diag}(\boldsymbol{\eta}(\boldsymbol{x})) - \boldsymbol{\eta}(\boldsymbol{x})\boldsymbol{\eta}(\boldsymbol{x})^{\mathsf{T}} \right) \mathbf{Y}^{\mathsf{T}} \qquad \forall \boldsymbol{x} \in \mathbb{R}^{d}.$$

$$(4.4)$$

On the one hand, since $\eta(x) \in \Delta_n$, (4.4) implies that

$$abla^2 arphi(oldsymbol{x}) \preceq rac{1}{ heta} \mathbf{Y} ext{diag}(oldsymbol{\eta}(oldsymbol{x})) \mathbf{Y}^{\mathsf{T}} \preceq rac{1}{ heta} \max_{1 \leq j \leq n} \left\{ \|oldsymbol{y}_j\|^2
ight\} \mathbf{I}_d \qquad orall oldsymbol{x} \in \mathbb{R}^d$$

On the other hand, for an arbitrary $\boldsymbol{x} \in \mathbb{R}^d$, let us define $\beta := n \min_{1 \le j \le n} \{\eta_j(\boldsymbol{x})\} \in (0, 1]$, let $\mathbf{1}_n$ denote the vector in \mathbb{R}^n with all entries equal to 1, let $\boldsymbol{p} := \frac{1}{1-\beta} (\boldsymbol{\eta}(\boldsymbol{x}) - \frac{\beta}{n} \mathbf{1}_n) \in \Delta_n$ if $\beta \neq 1$, and let $\boldsymbol{p} \in \Delta_n$

be arbitrary if $\beta = 1$. Then, it follows from the convexity of $\mathbb{R} \ni z \mapsto z^2 \in \mathbb{R}$ and Jensen's inequality that $\operatorname{diag}(\boldsymbol{p}) - \boldsymbol{p}\boldsymbol{p}^{\mathsf{T}} \succeq \mathbf{O}_n \mathbb{P}$ -almost surely. Observe that since $\boldsymbol{\eta}(\boldsymbol{x}) = \frac{\beta}{n} \mathbf{1}_n + (1 - \beta)\boldsymbol{p}$, it holds that

$$\begin{aligned} \operatorname{diag}(\boldsymbol{\eta}(\boldsymbol{x})) &- \boldsymbol{\eta}(\boldsymbol{x})\boldsymbol{\eta}(\boldsymbol{x})^{\mathsf{T}} \\ &= \frac{\beta}{n}\mathbf{I}_{n} + (1-\beta)\operatorname{diag}(\boldsymbol{p}) - \frac{\beta^{2}}{n^{2}}\mathbf{1}_{n}\mathbf{1}_{n}^{\mathsf{T}} - (1-\beta)^{2}\boldsymbol{p}\boldsymbol{p}^{\mathsf{T}} - \frac{\beta(1-\beta)}{n}\boldsymbol{p}\mathbf{1}_{n}^{\mathsf{T}} - \frac{\beta(1-\beta)}{n}\mathbf{1}_{n}\boldsymbol{p}^{\mathsf{T}} \\ &= \beta\left(\frac{1}{n}\mathbf{I}_{n} - \frac{1}{n^{2}}\mathbf{1}_{n}\mathbf{1}_{n}^{\mathsf{T}}\right) + (1-\beta)\left(\operatorname{diag}(\boldsymbol{p}) - \boldsymbol{p}\boldsymbol{p}^{\mathsf{T}}\right) + \beta(1-\beta)\left(\frac{1}{n}\mathbf{1}_{n} - \boldsymbol{p}\right)\left(\frac{1}{n}\mathbf{1}_{n} - \boldsymbol{p}\right)^{\mathsf{T}} \\ &\succeq \beta\left(\frac{1}{n}\mathbf{I}_{n} - \frac{1}{n^{2}}\mathbf{1}_{n}\mathbf{1}_{n}^{\mathsf{T}}\right). \end{aligned}$$

Subsequently, (4.4) implies that

$$\nabla^2 \varphi(\boldsymbol{x}) \succeq \frac{\beta}{\theta} \mathbf{Y} \left(\frac{1}{n} \mathbf{I}_n - \frac{1}{n^2} \mathbf{1}_n \mathbf{1}_n^{\mathsf{T}} \right) \mathbf{Y}^{\mathsf{T}} = \frac{\beta}{\theta} \widehat{\operatorname{Cov}}[\boldsymbol{y}_1, \dots, \boldsymbol{y}_n] \succeq \frac{\beta}{\theta} e_{\min} \big(\widehat{\operatorname{Cov}}[\boldsymbol{y}_1, \dots, \boldsymbol{y}_n] \big) \mathbf{I}_d.$$

The proof is now complete.

Proof of Proposition 4.1. Throughout this proof, let us denote

$$\Delta_{n} := \left\{ \boldsymbol{u} = (u_{1}, \dots, u_{n})^{\mathsf{T}} \in \mathbb{R}^{n} : \sum_{j=1}^{n} u_{j} = 1, \ u_{j} \ge 0 \ \forall 1 \le j \le n \right\},$$
$$\widehat{\eta}_{j}^{\theta}(\boldsymbol{x}) := \frac{\exp\left(\frac{1}{\theta}\left(\widehat{g}_{j}^{\theta} + \langle \boldsymbol{Y}_{j}, \boldsymbol{x} \rangle - \frac{1}{2} \|\boldsymbol{Y}_{j}\|^{2}\right)\right)}{\sum_{l=1}^{n} \exp\left(\frac{1}{\theta}\left(\widehat{g}_{l}^{\theta} + \langle \boldsymbol{Y}_{l}, \boldsymbol{x} \rangle - \frac{1}{2} \|\boldsymbol{Y}_{l}\|^{2}\right)\right)} \qquad \forall 1 \le j \le n, \ \forall \theta > 0, \ \forall \boldsymbol{x} \in \mathbb{R}^{d},$$
$$\widehat{\boldsymbol{\eta}}^{\theta}(\boldsymbol{x}) := \left(\widehat{\eta}_{1}^{\theta}(\boldsymbol{x}), \dots, \widehat{\eta}_{n}^{\theta}(\boldsymbol{x})\right)^{\mathsf{T}} \in \mathbb{R}^{n} \qquad \forall \theta > 0, \ \forall \boldsymbol{x} \in \mathbb{R}^{d}.$$

It holds that $\widehat{\eta}^{\theta}(\boldsymbol{x}) \in \Delta_n$ for all $\theta > 0$ and $\boldsymbol{x} \in \mathbb{R}^d$. Next, let us define $\widetilde{\varphi}_{entr}[\theta] : \mathbb{R}^d \to \mathbb{R}$ and $\widetilde{T}_{entr}[\theta] : \mathbb{R}^d \to \mathbb{R}^d$ as follows:

$$\begin{split} \widetilde{\varphi}_{\text{entr}}[\theta](\boldsymbol{x}) &:= \theta \log \left(\sum_{j=1}^{n} \exp\left(\frac{1}{\theta} \big(\widehat{g}_{j}^{\theta} + \langle \boldsymbol{Y}_{j}, \boldsymbol{x} \rangle - \frac{1}{2} \| \boldsymbol{Y}_{j} \|^{2} \big) \right) \right) \qquad \quad \forall \theta > 0, \ \forall \boldsymbol{x} \in \mathbb{R}^{d}, \\ \widetilde{T}_{\text{entr}}[\theta](\boldsymbol{x}) &:= \sum_{j=1}^{n} \widehat{\eta}_{j}^{\theta}(\boldsymbol{x}) \boldsymbol{Y}_{j} \qquad \qquad \quad \forall \theta > 0, \ \forall \boldsymbol{x} \in \mathbb{R}^{d}. \end{split}$$

Observe that $\widehat{\varphi}_{entr}[\theta](\boldsymbol{x}) = \widetilde{\varphi}_{entr}[\theta](\boldsymbol{x}) + \int_{0}^{\|\boldsymbol{x}\|^{2}/2} \gamma(z) dz$ and $\widehat{T}_{entr}[\theta](\boldsymbol{x}) = \widetilde{T}_{entr}[\theta](\boldsymbol{x}) + \gamma(\frac{1}{2}\|\boldsymbol{x}\|^{2})\boldsymbol{x}$ for all $\theta > 0$ and $\boldsymbol{x} \in \mathbb{R}^{d}$.

To begin, for any $\theta > 0$, Lemma 4.2 implies that $\widetilde{\varphi}_{entr}[\theta]$ is convex, $\widetilde{\varphi}_{entr}[\theta] \in \mathcal{C}^{\infty}(\mathbb{R}^d)$, and that $\widetilde{T}_{entr}[\theta] = \nabla \widetilde{\varphi}_{entr}[\theta]$. Moreover, the fundamental theorem of calculus implies that the gradient of $\mathbb{R}^d \ni \boldsymbol{x} \mapsto \int_0^{\|\boldsymbol{x}\|^2/2} \gamma(z) \, dz \in \mathbb{R}$ is $\mathbb{R}^d \ni \boldsymbol{x} \mapsto \gamma(\frac{1}{2}\|\boldsymbol{x}\|^2)\boldsymbol{x} \in \mathbb{R}^d$. Since $\gamma(\cdot)$ is non-negative, non-decreasing on \mathbb{R}_+ and infinitely differentiable on $(0, \infty)$, it holds that $\mathbb{R}^d \ni \boldsymbol{x} \mapsto \int_0^{\|\boldsymbol{x}\|^2/2} \gamma(z) \, dz \in \mathbb{R}$ is convex and infinitely differentiable on \mathbb{R}^d . Therefore, we get $\widehat{\varphi}_{entr}[\theta]$ is convex and infinitely differentiable, and $\widehat{T}_{entr}[\theta] = \nabla \widehat{\varphi}_{entr}[\theta]$. Since $\gamma(\cdot)$ is bounded and Lemma 4.2 shows that $\nabla \widetilde{T}_{entr}[\theta] = \nabla^2 \widetilde{\varphi}_{entr}[\theta] \preceq \frac{1}{\theta} \max_{1 \le j \le n} \{\|\boldsymbol{Y}_j\|^2\} \mathbf{I}_d$, we get $\widehat{T}_{entr}[\theta] \in \mathcal{C}_{lin}(\mathbb{R}^d, \mathbb{R}^d)$. Moreover, [56, Example 14.31 & Theorem 14.37] implies that $(\widehat{g}_j^\theta)_{j=1:n}$ have a Borel dependence on $(\boldsymbol{X}_1, \ldots, \boldsymbol{X}_m, \boldsymbol{Y}_1, \ldots, \boldsymbol{Y}_n, \theta)$. Consequently, $\widehat{T}_{entr}[\boldsymbol{X}_1, \ldots, \boldsymbol{X}_m, \boldsymbol{Y}_1, \ldots, \boldsymbol{Y}_n, \theta]$.

Let us fix arbitrary $m \ge 1$, $n \ge d + 1$, $\theta > 0$ and prove statement (i). Since we have shown that $\widehat{\varphi}_{entr}[\theta]$ is convex and infinitely differentiable, and $\widehat{T}_{entr}[\theta] = \nabla \widehat{\varphi}_{entr}[\theta]$, it remains to show that $\widehat{\varphi}_{entr}[\theta] \in \mathfrak{C}_{\underline{\lambda},\infty}(\mathbb{R}^d)$ and $\underline{\lambda} > 0$ \mathbb{P} -almost surely. To that end, we begin by fixing an arbitrary $x \in \overline{B}(\mathbf{0}, 2r_0(\mu))$ and observing from the first-order optimality conditions of the dual entropic optimal transport problem (4.1) that $\frac{1}{m} \sum_{i=1}^{m} \exp\left(\frac{1}{\theta}(\widehat{f}_i^{\theta} + 1) + 1\right)$

 $\widehat{g}_j^{\theta} - \frac{1}{2} \| \boldsymbol{X}_i - \boldsymbol{Y}_j \|^2 \Big) = 1$ for $j = 1, \dots, n$, and hence

$$\exp\left(\frac{1}{\theta}\widehat{g}_{j}^{\theta}\right) = \left(\frac{1}{m}\sum_{i=1}^{m}\exp\left(\frac{1}{\theta}\left(\widehat{f}_{i}^{\theta} - \frac{1}{2}\|\boldsymbol{X}_{i} - \boldsymbol{Y}_{j}\|^{2}\right)\right)\right)^{-1} \quad \forall 1 \le j \le n.$$
(4.5)

Since it holds that

$$\begin{split} \left| \| \boldsymbol{X}_{i} - \boldsymbol{Y}_{j} \|^{2} - \| \boldsymbol{X}_{i} - \boldsymbol{Y}_{l} \|^{2} \right| &= \left(\| \boldsymbol{X}_{i} - \boldsymbol{Y}_{j} \| + \| \boldsymbol{X}_{i} - \boldsymbol{Y}_{l} \| \right) \left| \| \boldsymbol{X}_{i} - \boldsymbol{Y}_{j} \| - \| \boldsymbol{X}_{i} - \boldsymbol{Y}_{l} \| \right| \\ &\leq \left(2 \| \boldsymbol{X}_{i} \| + \| \boldsymbol{Y}_{j} \| + \| \boldsymbol{Y}_{l} \| \right) \left(\| \boldsymbol{Y}_{j} \| + \| \boldsymbol{Y}_{l} \| \right) \\ &\leq \left(4 r_{0}(\mu) + 4 r_{0}(\nu) \right) r_{0}(\nu) \qquad \mathbb{P}\text{-a.s. } \forall 1 \leq i \leq m, \ \forall 1 \leq j \leq n, \ \forall 1 \leq l \leq n, \end{split}$$

we get from (4.5) that

$$\exp\left(\frac{1}{\theta}\widehat{g}_{j}^{\theta}\right) \geq \left(\frac{1}{m}\sum_{i=1}^{m}\exp\left(\frac{1}{\theta}\left(\widehat{f}_{i}^{\theta}-\frac{1}{2}\|\boldsymbol{X}_{i}-\boldsymbol{Y}_{l}\|^{2}+\frac{1}{2}\|\|\boldsymbol{X}_{i}-\boldsymbol{Y}_{j}\|^{2}-\|\boldsymbol{X}_{i}-\boldsymbol{Y}_{l}\|^{2}|\right)\right)\right)^{-1}$$

$$\geq \left(\frac{1}{m}\sum_{i=1}^{m}\exp\left(\frac{1}{\theta}\left(\widehat{f}_{i}^{\theta}-\frac{1}{2}\|\boldsymbol{X}_{i}-\boldsymbol{Y}_{l}\|^{2}\right)\right)\right)^{-1}\exp\left(-\frac{1}{\theta}\left(2r_{0}(\mu)+2r_{0}(\nu)\right)r_{0}(\nu)\right)$$

$$\geq \exp\left(\frac{1}{\theta}\widehat{g}_{l}^{\theta}\right)\exp\left(-\frac{1}{\theta}\left(2r_{0}(\mu)+2r_{0}(\nu)\right)r_{0}(\nu)\right) \qquad \mathbb{P}\text{-a.s. } \forall 1 \leq j \leq n, \ \forall 1 \leq l \leq n.$$

$$(4.6)$$

Moreover, since $\|\boldsymbol{x}\| \leq 2r_0(\mu)$, it holds that

$$\begin{split} \left| \| \boldsymbol{x} - \boldsymbol{Y}_{j} \|^{2} - \| \boldsymbol{x} - \boldsymbol{Y}_{l} \|^{2} \right| &= \left(\| \boldsymbol{x} - \boldsymbol{Y}_{j} \| + \| \boldsymbol{x} - \boldsymbol{Y}_{l} \| \right) \left| \| \boldsymbol{x} - \boldsymbol{Y}_{j} \| - \| \boldsymbol{x} - \boldsymbol{Y}_{l} \| \right| \\ &\leq \left(2 \| \boldsymbol{x} \| + \| \boldsymbol{Y}_{j} \| + \| \boldsymbol{Y}_{l} \| \right) \left(\| \boldsymbol{Y}_{j} \| + \| \boldsymbol{Y}_{l} \| \right) \\ &\leq \left(8 r_{0}(\mu) + 4 r_{0}(\nu) \right) r_{0}(\nu) \qquad \mathbb{P}\text{-a.s. } \forall 1 \leq j \leq n, \ \forall 1 \leq l \leq n, \end{split}$$

which then yields

$$\exp\left(-\frac{1}{2\theta}\|\boldsymbol{x}-\boldsymbol{Y}_{j}\|^{2}\right) \geq \exp\left(-\frac{1}{2\theta}\left(\|\boldsymbol{x}-\boldsymbol{Y}_{l}\|^{2}+\left\|\|\boldsymbol{x}-\boldsymbol{Y}_{j}\|^{2}-\|\boldsymbol{x}-\boldsymbol{Y}_{l}\|^{2}\right)\right)\right)$$
$$\geq \exp\left(-\frac{1}{2\theta}\|\boldsymbol{x}-\boldsymbol{Y}_{l}\|^{2}\right)\exp\left(-\frac{1}{\theta}\left(4r_{0}(\mu)+2r_{0}(\nu)\right)r_{0}(\nu)\right)$$
$$\mathbb{P}\text{-a.s. }\forall 1 \leq j \leq n, \ \forall 1 \leq l \leq n.$$

$$(4.7)$$

Combining (4.6) and (4.7) leads to

$$\exp\left(\frac{1}{\theta}\left(\widehat{g}_{j}^{\theta}-\frac{1}{2}\|\boldsymbol{x}-\boldsymbol{Y}_{j}\|^{2}\right)\right) \geq \exp\left(\frac{1}{\theta}\left(\widehat{g}_{l}^{\theta}-\frac{1}{2}\|\boldsymbol{x}-\boldsymbol{Y}_{l}\|^{2}\right)\right)\exp\left(-\frac{1}{\theta}\left(6r_{0}(\mu)+4r_{0}(\nu)\right)r_{0}(\nu)\right)$$
$$\mathbb{P}\text{-a.s. }\forall 1 \leq j \leq n, \ \forall 1 \leq l \leq n.$$

Consequently, we get

$$\begin{split} \widehat{\eta}_{j}^{\theta}(\boldsymbol{x}) &= \frac{\exp\left(\frac{1}{\theta}\left(\widehat{g}_{j}^{\theta} + \langle \boldsymbol{Y}_{j}, \boldsymbol{x} \rangle - \frac{1}{2} \|\boldsymbol{Y}_{j}\|^{2}\right)\right)}{\sum_{l=1}^{n} \exp\left(\frac{1}{\theta}\left(\widehat{g}_{l}^{\theta} + \langle \boldsymbol{Y}_{l}, \boldsymbol{x} \rangle - \frac{1}{2} \|\boldsymbol{Y}_{l}\|^{2}\right)\right)} \\ &= \frac{\exp\left(\frac{1}{\theta}\left(\widehat{g}_{j}^{\theta} - \frac{1}{2} \|\boldsymbol{x} - \boldsymbol{Y}_{j}\|^{2}\right)\right)}{\sum_{l=1}^{n} \exp\left(\frac{1}{\theta}\left(\widehat{g}_{l}^{\theta} - \frac{1}{2} \|\boldsymbol{x} - \boldsymbol{Y}_{l}\|^{2}\right)\right)} \\ &\geq \frac{1}{n} \exp\left(-\frac{1}{\theta}\left(6r_{0}(\mu) + 4r_{0}(\nu)\right)r_{0}(\nu)\right) \quad \mathbb{P}\text{-a.s. } \forall 1 \leq j \leq n \end{split}$$

Now, letting $\beta := \frac{1}{n} \exp\left(-\frac{1}{\theta} (6r_0(\mu) + 4r_0(\nu))r_0(\nu)\right) \le \min_{1 \le j \le n} \left\{\widehat{\eta}_j^{\theta}(\boldsymbol{x})\right\}$, we get from Lemma 4.2 that $\nabla^2 \widetilde{\varphi}_{entr}[\theta](\boldsymbol{x}) \succeq \frac{n\beta}{\theta} e_{\min}(\widehat{Cov}[\boldsymbol{Y}_1, \dots, \boldsymbol{Y}_n])\mathbf{I}_d \succeq \underline{\lambda}\mathbf{I}_d \qquad \mathbb{P}\text{-a.s.}$ Thus, we have shown that $\nabla^2 \widehat{\varphi}_{entr}[\theta](\boldsymbol{x}) \succeq \nabla^2 \widetilde{\varphi}_{entr}[\theta](\boldsymbol{x}) \succeq \underline{\lambda} \mathbf{I}_d$ for all $\boldsymbol{x} \in \overline{B}(\boldsymbol{0}, 2r_0(\mu))$ \mathbb{P} -almost surely. Lastly, for any $\boldsymbol{x} \in \mathbb{R}^d \setminus \overline{B}(\boldsymbol{0}, 2r_0(\mu))$, since $\nabla^2 \widehat{\varphi}_{entr}[\theta](\boldsymbol{x}) \succeq \mathbf{O}_d$ and $\gamma(\frac{1}{2} \|\boldsymbol{x}\|^2) \ge \gamma(2r_0(\mu)) = \exp\left(-\frac{2}{3r_0(\mu)^2}\right)$, it holds that

$$\nabla^2 \widehat{\varphi}_{entr}[\theta](\boldsymbol{x}) = \nabla^2 \widetilde{\varphi}_{entr}[\theta](\boldsymbol{x}) + \gamma' (\frac{1}{2} \|\boldsymbol{x}\|^2) \boldsymbol{x} \boldsymbol{x}^{\mathsf{T}} + \gamma (\frac{1}{2} \|\boldsymbol{x}\|^2) \mathbf{I}_d$$
$$\succeq \gamma (\frac{1}{2} \|\boldsymbol{x}\|^2) \mathbf{I}_d \succeq \exp\left(-\frac{2}{3r_0(\mu)^2}\right) \mathbf{I}_d \succeq \underline{\lambda} \mathbf{I}_d.$$

Since ν is absolutely continuous respect to the Lebesgue measure on \mathbb{R}^d and $n \ge d+1$, the sample covariance matrix $\widehat{\text{Cov}}[\mathbf{Y}_1, \ldots, \mathbf{Y}_n]$ is \mathbb{P} -almost surely non-singular. Therefore, it holds \mathbb{P} -almost surely that $\underline{\lambda}(\mu, \nu, m, n, \mathbf{X}_1, \ldots, \mathbf{X}_m, \mathbf{Y}_1, \ldots, \mathbf{Y}_n, \theta) > 0$, and the proof of statement (i) is now complete.

To prove statement (ii), let us fix arbitrary $m, n \in \mathbb{N}, \theta > 0$, as well as an arbitrary $x \in \mathbb{R}^d$. Observe that

$$egin{aligned} &\|\widehat{T}_{ ext{entr}}[heta](m{x}) - \widehat{T}_{ ext{entr}}[heta](m{0})\|^2 &\leq 2 \|\widetilde{T}_{ ext{entr}}[heta](m{x}) - \widetilde{T}_{ ext{entr}}[heta](m{0})\|^2 + 2\gamma ig(rac{1}{2} \|m{x}\|^2ig)\|m{x}\|^2 &\leq 2 \|\widetilde{T}_{ ext{entr}}[heta](m{x}) - \widetilde{T}_{ ext{entr}}[heta](m{0})\|^2 + 2\|m{x}\|^2. \end{aligned}$$

Since both $\widetilde{T}_{entr}[\theta](\boldsymbol{x})$ and $\widetilde{T}_{entr}[\theta](\boldsymbol{0})$ belong to $\operatorname{conv}(\{\boldsymbol{Y}_1,\ldots,\boldsymbol{Y}_n\})$ by definition, it holds \mathbb{P} -almost surely that $\|\widetilde{T}_{entr}[\theta](\boldsymbol{x}) - \widetilde{T}_{entr}[\theta](\boldsymbol{0})\| \leq 2 \sup_{\boldsymbol{y} \in \operatorname{supp}(\nu)} \{\|\boldsymbol{y}\|\} \leq 2r_0(\nu)$. This shows that $\|\widehat{T}_{entr}[\theta](\boldsymbol{x}) - \widehat{T}_{entr}[\theta](\boldsymbol{0})\|^2 \leq 8r_0(\nu)^2 + 2\|\boldsymbol{x}\|^2$ for all $\boldsymbol{x} \in \mathbb{R}^d$ \mathbb{P} -almost surely, which then yields $\mathbb{E}[\|\widehat{T}_{entr}[\theta](\boldsymbol{x}) - \widehat{T}_{entr}[\theta](\boldsymbol{0})\|^2] \leq 8r_0(\nu)^2 + 2\|\boldsymbol{x}\|^2$ for all $\boldsymbol{x} \in \mathbb{R}^d$. This completes the proof of statement (ii).

In the following, we will prove statement (iii). First, observe that $\gamma(\frac{1}{2} || \boldsymbol{x} ||^2) = 0$ for all $\boldsymbol{x} \in \bar{B}(\boldsymbol{0}, r_0(\mu))$. In particular, $\gamma(\frac{1}{2} || \boldsymbol{x} ||^2) = 0$ for all $\boldsymbol{x} \in \operatorname{supp}(\mu)$, and hence

$$\mathbb{E}\Big[\big\|\widehat{T}_{\text{entr}}[\theta] - T^{\mu}_{\nu}\big\|_{\mathcal{L}^{2}(\mu)}^{2}\Big] = \mathbb{E}\Big[\big\|\widetilde{T}_{\text{entr}}[\theta] - T^{\mu}_{\nu}\big\|_{\mathcal{L}^{2}(\mu)}^{2}\Big] \qquad \forall m, n \in \mathbb{N}, \ \forall \theta > 0.$$

Next, we have by the assumption that $\mu, \nu \in \mathcal{M}^q(\mathbb{R}^d)$ and by Caffarelli's regularity theory (Theorem 2.5) that the Brenier potentials φ_{ν}^{μ} and φ_{ν}^{ν} satisfy $\varphi_{\nu}^{\mu} \in C^{q+2,\alpha}(\operatorname{supp}(\mu))$ and $\varphi_{\mu}^{\nu} \in C^{q+2,\alpha}(\operatorname{supp}(\nu))$ for some $\alpha \in (0, 1]$. Let us define $\overline{\alpha}(\mu, \nu) := \min\{q + 2 + \alpha, 4\} \in [3, 4]$. In addition, Lemma 3.2 implies that there exist $0 < \lambda_{\text{LB}} \leq \lambda_{\text{UB}} < \infty$ such that $\lambda_{\text{LB}}\mathbf{I}_d \preceq \nabla^2 \varphi_{\nu}^{\mu}(\mathbf{x}) \preceq \lambda_{\text{UB}}\mathbf{I}_d$ for all $\mathbf{x} \in \operatorname{supp}(\mu)$. Thus, one may check the assumptions (A1)–(A3) in [52] are satisfied with respect to μ and ν . It subsequently follows from [52, Theorem 4 & Theorem 5] that there exist $C_1(\mu, \nu) > 0$, $C_2(\mu, \nu) > 0$, $C_3(\mu, \nu) > 0$, $C_4(\mu, \nu) > 0$, $\overline{\theta}(\mu, \nu) > 0$ that only depend on μ and ν such that

$$\mathbb{E}\Big[\|\widetilde{T}_{entr}[\theta] - T^{\mu}_{\nu}\|^{2}_{\mathcal{L}^{2}(\mu)} \Big] \leq C_{1}(\mu,\nu)\theta^{1-\frac{d}{2}}\log(n)n^{-\frac{1}{2}} + C_{2}(\mu,\nu)\theta^{\frac{\overline{\alpha}(\mu,\nu)}{2}}
+ C_{3}(\mu,\nu)\theta^{2}I_{0}(\mu,\nu) + C_{4}(\mu,\nu)\theta^{-\frac{d}{2}}\log(m)m^{-\frac{1}{2}}
\forall m \geq 2, \ \forall n \geq 2, \ \forall 0 < \theta \leq \overline{\theta}(\mu,\nu),$$
(4.8)

where $I_0(\mu, \nu)$ is the integrated Fisher information along the Wasserstein geodesic between μ and ν defined in Appendix A of [52]. Since $\varphi_{\nu}^{\mu} \in C^{q+2,\alpha}(\operatorname{supp}(\mu))$ where $q+2 \geq 3$, $\nabla^2 \varphi_{\nu}^{\mu}$ is Lipschitz continuous, and we have by [15, Proposition 1] that $I_0(\mu, \nu) < \infty$. The above bound can thus be further bounded as follows:

$$C_{1}(\mu,\nu)\theta^{1-\frac{d}{2}}\log(n)n^{-\frac{1}{2}} + C_{2}(\mu,\nu)\theta^{\frac{\overline{\alpha}(\mu,\nu)}{2}} + C_{3}(\mu,\nu)\theta^{2}I_{0}(\mu,\nu) + C_{4}(\mu,\nu)\theta^{-\frac{d}{2}}\log(m)m^{-\frac{1}{2}}$$

$$\leq C_{1}(\mu,\nu)\overline{\theta}(\mu,\nu)\theta^{-\frac{d}{2}}\log(n)n^{-\frac{1}{2}} + C_{4}(\mu,\nu)\theta^{-\frac{d}{2}}\log(m)m^{-\frac{1}{2}}$$

$$+ \left(C_{2}(\mu,\nu) + C_{3}(\mu,\nu)\overline{\theta}(\mu,\nu)^{\frac{4-\overline{\alpha}(\mu,\nu)}{2}}I_{0}(\mu,\nu)\right)\theta^{\frac{\overline{\alpha}(\mu,\nu)}{2}}$$

$$\leq C_{\text{entr}}(\mu,\nu)\left[\theta^{-\frac{d}{2}}\left(\log(\min\{m,n\})\min\{m,n\}^{-\frac{1}{2}}\right) + \theta^{\frac{\overline{\alpha}(\mu,\nu)}{2}}\right]$$

$$\forall m \geq 7, \ \forall n \geq 7, \ \forall 0 < \theta \leq \overline{\theta}(\mu,\nu),$$
(4.9)

where $C_{\text{entr}}(\mu,\nu) := \max\left\{C_1(\mu,\nu)\overline{\theta}(\mu,\nu) + C_4(\mu,\nu), C_2(\mu,\nu) + C_3(\mu,\nu)\overline{\theta}(\mu,\nu)^{\frac{4-\overline{\alpha}(\mu,\nu)}{2}}I_0(\mu,\nu)\right\} < \infty$. Now, let us fix arbitrary $\epsilon > 0, m \ge \overline{n}(\mu,\nu,\epsilon)$, and $n \ge \overline{n}(\mu,\nu,\epsilon)$. Recall that the definition of $\overline{n}(\mu,\nu,\epsilon)$ guarantees $m \ge 7, n \ge 7, n \ge d+1$, and $\min\{m,n\}^{-\frac{\overline{\alpha}(\mu,\nu)}{2(\overline{\alpha}(\mu,\nu)+d)}} \left(\log\left(\min\{m,n\}\right)+1\right) \le \frac{\epsilon}{C_{\text{entr}}(\mu,\nu)}$. We thus get $0 < \tilde{\theta}(\mu,\nu,m,n,\epsilon) = \min\{m,n\}^{-\frac{1}{\overline{\alpha}(\mu,\nu)+d}} \le \overline{n}(\mu,\nu,\epsilon)^{-\frac{1}{\overline{\alpha}(\mu,\nu)+d}} \le \overline{\theta}(\mu,\nu)$. Subsequently, combining (4.8) and (4.9) yields

$$\mathbb{E}\Big[\big\|\widetilde{T}_{\text{entr}}\big[\widetilde{\theta}(\mu,\nu,m,n,\epsilon)\big] - T^{\mu}_{\nu}\big\|_{\mathcal{L}^{2}(\mu)}^{2}\Big]$$

$$\leq C_{\text{entr}}(\mu,\nu)\Big(\widetilde{\theta}(\mu,\nu,m,n,\epsilon)^{-\frac{d}{2}}\log\big(\min\{m,n\}\big)\min\{m,n\}^{-\frac{1}{2}} + \widetilde{\theta}(\mu,\nu,m,n,\epsilon)^{\frac{\overline{\alpha}(\mu,\nu)}{2}}\Big)$$

$$= C_{\text{entr}}(\mu,\nu)\Big(\log\big(\min\{m,n\}\big)\min\{m,n\}^{-\frac{1}{2}+\frac{d}{2(\overline{\alpha}(\mu,\nu)+d)}} + \min\{m,n\}^{-\frac{\overline{\alpha}(\mu,\nu)}{2(\overline{\alpha}(\mu,\nu)+d)}}\Big)$$

$$\leq C_{\text{entr}}(\mu,\nu)\Big(\log\big(\min\{m,n\}\big) + 1\Big)\min\{m,n\}^{-\frac{\overline{\alpha}(\mu,\nu)}{2(\overline{\alpha}(\mu,\nu)+d)}} \leq \epsilon.$$

The proof is now complete.

Since Proposition 4.1 has shown that $\widehat{T}_{entr}[\theta]$ satisfies Assumption 3.4, letting $\widehat{T}_{\nu,n}^{\mu,m}[\theta] \leftarrow \widehat{T}_{entr}[\theta]$ in Algorithm 2 leads to the \mathbb{P} -almost sure convergence of the output $(\widehat{\mu}_t)_{t\in\mathbb{N}_0}$ under Setting 3.13. This is summarized in the following corollary.

Corollary 4.3 (Convergence of Algorithm 2 with $\widehat{T}_{entr}[\theta]$). Let the inputs of Algorithm 2 satisfy Setting 3.6, let the OT map estimator $\widehat{T}_{\nu,n}^{\mu,m}[\theta]$ be given by $\widehat{T}_{entr}[\theta]$ defined in Proposition 4.1, and let $u_0(\nu)$, $u_1(\nu)$, $\overline{n}(\mu,\nu,\epsilon)$, $\widetilde{\theta}(\mu,\nu,m,n,\epsilon)$ be defined as in Proposition 4.1. Moreover, let $(\Omega, \mathcal{F}, \mathbb{P}, (\mathcal{F}_t)_{t\in\mathbb{N}_0})$ be the filtered probability space constructed by Algorithm 2, let the $(\mathcal{F}_t)_{t\in\mathbb{N}_0}$ -adapted stochastic processes $(\widehat{R}_t)_{t\in\mathbb{N}_0}$, $(\widehat{N}_{t,k})_{k=1:K, t\in\mathbb{N}_0}$, and $(\widehat{\Theta}_{t,k})_{k=1:K, t\in\mathbb{N}_0}$ in Algorithm 2 be specified by Setting 3.13, and let $(\widehat{\mu}_t)_{t\in\mathbb{N}_0}$ be the output of Algorithm 2. Then, the following statements hold.

- (i) It holds \mathbb{P} -almost surely that $(\widehat{\mu}_t)_{t\in\mathbb{N}_0}$ is precompact with respect to the \mathcal{W}_2 -metric. Moreover, every accumulation point of $(\widehat{\mu}_t)_{t\in\mathbb{N}_0}$ with respect to the \mathcal{W}_2 -metric is a fixed-point of G.
- (ii) In particular, if G has a unique fixed-point, then $(\widehat{\mu}_t)_{t\in\mathbb{N}_0}$ converges \mathbb{P} -almost surely in \mathcal{W}_2 to the Wasserstein barycenter of ν_1, \ldots, ν_K with weights w_1, \ldots, w_K .

Remark 4.4 (Computational tractability of $\widehat{T}_{entr}[\theta]$). The computation of $\widehat{T}_{entr}[\theta]$ is done in two phases. In the first phase, given the m samples X_1, \ldots, X_m from μ and the n samples Y_1, \ldots, Y_n from ν , one computes an optimal solution $(\widehat{f}_i^{\theta})_{i=1:m}, (\widehat{g}_j^{\theta})_{j=1:n}$ of the dual entropic optimal transport problem (4.1) by the well-known Sinkhorn's algorithm [20]. Sinkhorn's algorithm is scalable to large sample sizes and admits highly parallelizable implementations on GPUs; see, e.g., [51, Section 4.3] for the numerical aspects of Sinkhorn's algorithm, and [27] for a multiscale extension. In our numerical experiments, we utilize the Optimal Transport Tools (OTT) Python toolbox developed by Cuturi, Meng-Papaxanthos, Tian, Bunne, Davis, and Teboul [21] to solve (4.1). In the second phase, one can efficiently evaluate $\widehat{T}_{entr}[\theta](x)$ at any $x \in \mathbb{R}^d$ directly through its definition in (4.3). Consequently, using $\widehat{T}_{entr}[\theta]$ as the OT map estimator in Algorithm 2 results in a computationally tractable and efficient algorithm for Wasserstein barycenter that is also provably convergent.

Remark 4.5 (An alternative choice of admissible OT map estimator). Apart from the entropic OT map estimator introduced in this section, one could also adopt the convex least squares estimator of Manole et al. [42, Proposition 16], although it needs to be appropriately modified to possess the strong-convexity and differentiability properties in Assumption 3.4(i) (e.g., by imposing shape constraints). However, in this paper, we choose to focus on the entropic OT map estimator due to its superior computational efficiency.

5. A NOVEL ALGORITHM FOR GENERATING SYNTHETIC PROBLEM INSTANCES

Empirical experiments on computing the W_2 -barycenter for continuous measures in most existing studies have been restricted to problem instances where the input measures belong to the same family of elliptical distributions (see, e.g., [43, Definition 3.26] for the definition), of which the ground-truth W_2 -barycenter can be accurately and efficiently computed. To evaluate approximated W_2 -barycenters for non-elliptical measures, a common practice is to conduct experiments on low-dimensional imaging datasets and visually assess the image generated from the approximated W_2 -barycenters. However, such approaches purely rely on human judgement and lack explicit numerical evidence without access to the ground-truth barycenters. Therefore, it is important to develop problem instances with non-parametric free-support input measures where the groundtruth W_2 -barycenter is a priori known, such that quantitative inspections of empirical approximation errors can be conducted.

Algorithm 3: Synthetic generation of Wasserstein barycenter problem instance.

 $\overline{\text{Input: } \bar{\mu} \in \mathcal{P}_{2, \mathrm{ac}}(\mathbb{R}^d), K \in \mathbb{N}, w_1 > 0, \dots, w_K > 0, \widetilde{K} \in \mathbb{N}, (\underline{\lambda}_{\widetilde{k}}, \theta_{\widetilde{k}}, \alpha_{\widetilde{k}}, n_{\widetilde{k}}, (g_{\widetilde{k}, j}, \boldsymbol{y}_{\widetilde{k}, j})_{j=1:n_{\widetilde{k}}})_{\widetilde{k}=1:\widetilde{K}}, }$ $\Phi: \{1,\ldots,\widetilde{K},-1,\ldots,-\widetilde{K}\} \to \{1,\ldots,K\}, (\mathbf{A}_k, \mathbf{b}_k)_{k=1:K}, \gamma \in [0,1),$ TRUNCATE \in {True, False}, and $(\mathcal{Y}_k)_{k=1:K} \subset \mathcal{S}(\mathbb{R}^d)$ satisfying Setting 5.1. **Output:** $(\nu_k, T_k)_{k=1:K}$. 1 for k = 1, ..., K do For $j = 1, \ldots, n_{\widetilde{k}}$, define $\eta_{\widetilde{k}, j}(\boldsymbol{x}) := \frac{\exp\left(\frac{1}{\theta_{\widetilde{k}}}\left(g_{\widetilde{k}, j} + \langle \boldsymbol{y}_{\widetilde{k}, j}, \boldsymbol{x} \rangle - \frac{1}{2} \|\boldsymbol{y}_{\widetilde{k}, j}\|^2\right)\right)}{\sum_{j'=1}^{n_{\widetilde{k}}} \exp\left(\frac{1}{\theta_{\widetilde{k}}}\left(g_{\widetilde{k}, j'} + \langle \boldsymbol{y}_{\widetilde{k}, j'}, \boldsymbol{x} \rangle - \frac{1}{2} \|\boldsymbol{y}_{\widetilde{k}, j'}\|^2\right)\right)}$ for all $\boldsymbol{x} \in \mathbb{R}^d$. 2 Define $\boldsymbol{\eta}_{\widetilde{k}}(\boldsymbol{x}) := \left(\eta_{\widetilde{k},1}(\boldsymbol{x}), \ldots, \eta_{\widetilde{k},n_{\widetilde{k}}}(\boldsymbol{x})\right)^{\mathsf{T}} \in \mathbb{R}^{n_{\widetilde{k}}^{\times}}$ for all $\boldsymbol{x} \in \mathbb{R}^{d}$ 3 Define $\mathbf{Y}_{\widetilde{k}} := \left(\begin{array}{ccc} | & | & | \\ y_{\widetilde{k},1} & y_{\widetilde{k},2} & \cdots & y_{\widetilde{k},n_{\widetilde{k}}} \end{array}
ight) \in \mathbb{R}^{d \times n_{\widetilde{k}}}.$ 4 Choose $\overline{\lambda}_{\widetilde{k}} \geq \frac{1}{\theta_{\widetilde{k}}} \max_{\boldsymbol{x} \in \mathbb{R}^d} \left\{ e_{\max} \left(\mathbf{Y}_{\widetilde{k}} (\operatorname{diag}(\boldsymbol{\eta}_{\widetilde{k}}(\boldsymbol{x})) - \boldsymbol{\eta}_{\widetilde{k}}(\boldsymbol{x}) \boldsymbol{\eta}_{\widetilde{k}}(\boldsymbol{x})^{\mathsf{T}} \right) \mathbf{Y}_{\widetilde{k}}^{\mathsf{T}} \right) \right\} + 2\underline{\lambda}_{\widetilde{k}}.$ 5 Define $\widetilde{T}_{\widetilde{k}}(\boldsymbol{x}) := \left(\sum_{j=1}^{n_{\widetilde{k}}} \eta_{\widetilde{k},j}(\boldsymbol{x})\boldsymbol{y}_{\widetilde{k},j}\right) + \underline{\lambda}_{\widetilde{k}}\boldsymbol{x} \text{ and } \widetilde{T}_{-\widetilde{k}}(\boldsymbol{x}) := \overline{\lambda}_{\widetilde{k}}\boldsymbol{x} - \widetilde{T}_{\widetilde{k}}(\boldsymbol{x}) \text{ for all } \boldsymbol{x} \in \mathbb{R}^d.$ 6 7 for $\widetilde{k} = 1, \ldots, \widetilde{K}$ do Set $\beta_{-\widetilde{k}} \leftarrow (1-\gamma) \alpha_{\widetilde{k}} \Big(\sum_{\widetilde{k}'=1}^{\widetilde{K}} w_{\Phi(-\widetilde{k}')} \alpha_{\widetilde{k}'} \overline{\lambda}_{\widetilde{k}'} \Big)^{-1} \in (0,\infty), \beta_{\widetilde{k}} \leftarrow \frac{w_{\Phi(-\widetilde{k})}}{w_{\Phi(\widetilde{k})}} \beta_{-\widetilde{k}} \in (0,\infty).$ 8 9 for $k = 1, \ldots, K$ do Define $T_k(\boldsymbol{x}) := \left(\sum_{i \in \Phi^{-1}(k)} \beta_i \widetilde{T}_i(\boldsymbol{x})\right) + \gamma(\mathbf{A}_k \boldsymbol{x} + \boldsymbol{b}_k)$ for all $\boldsymbol{x} \in \mathbb{R}^d$. 10 if TRUNCATE = True then 11 Set $\nu_k \leftarrow (T_k \sharp \bar{\mu}) |_{\mathcal{Y}_k}$. 12 else 13 Set $\nu_k \leftarrow T_k \sharp \bar{\mu}$. 14 15 return $(\nu_k, T_k)_{k=1:K}$.

To this end, Korotin et al. [38] proposed a method of generating input measures using an initial measure which ends up being exactly the W_2 -barycenter, via exploiting the convexity and congruency properties inherited by the Brenier potential functions (see [38, Section 5]). Although their method serves as a reasonable benchmark in many computer vision and imaging applications (e.g., color transfer), the conjugacy operation therein can be computationally challenging, and the constructed congruent functions suffer from limited curvatures. As a consequence, the resulting input measures exhibit little structural differences between each other and are close to the pushforwards of the initial measure under certain close-to-affine transformations, which hinders the generalizability of the problem instance.

In this section, we present a novel and flexible algorithm for synthetically generating problem instances that can be used for evaluating the efficacy of Wasserstein barycenter algorithms. In particular, our method is inspired by the method of Korotin et al. [38] to generate a problem instance out of a known measure $\bar{\mu}$ as the underlying Wasserstein barycenter. However, it is computationally more efficient, allows for arbitrary barycenter weights $w_1 > 0, \ldots, w_K > 0$ that can be user-specified, and creates input measures ν_1, \ldots, ν_K with more non-trivial structures.

Specifically, the input measures ν_1, \ldots, ν_K of our generated problem instances will be characterized via the pushforwards of a known probability measure $\bar{\mu} \in \mathcal{P}_{2,ac}(\mathbb{R}^d)$ by several tailored transportation maps $T_1, \ldots, T_K : \mathbb{R}^d \to \mathbb{R}^d$. The generation procedure is detailed in Algorithm 3, whose inputs are specified in the following setting.

Setting 5.1 (Inputs of Algorithm 3). Let $\bar{\mu} \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$, let $K \in \mathbb{N} \cap [2, \infty)$, and let $w_1 > 0, \ldots, w_K > 0$ satisfy $\sum_{k=1}^{K} w_k = 1$. Let $\tilde{K} \in \mathbb{N} \cap [2, \infty)$ satisfy $2\tilde{K} \ge K$, and for $\tilde{k} = 1, \ldots, \tilde{K}$, let $\underline{\lambda}_{\tilde{k}} > 0$, $\theta_{\tilde{k}} > 0$, $\alpha_{\tilde{k}} > 0$, $n_{\tilde{k}} \in \mathbb{N}$, $(g_{\tilde{k},j})_{j=1:n_{\tilde{k}}} \subset \mathbb{R}$, and $(\boldsymbol{y}_{\tilde{k},j})_{j=1:n_{\tilde{k}}} \subset \mathbb{R}^d$. Moreover, let $\Phi : \{1, \ldots, \tilde{K}, -1, \ldots, -\tilde{K}\} \rightarrow \{1, \ldots, K\}$ be a surjective map. Furthermore, let $(\mathbf{A}_k)_{k=1:K} \subset \mathbb{S}^d_{++}, (\boldsymbol{b}_k)_{k=1:K} \subset \mathbb{R}^d$ satisfy $\sum_{k=1}^{K} w_k \mathbf{A}_k = \mathbf{I}_d$ and $\sum_{k=1}^{K} w_k b_k = \mathbf{0}_d$. Lastly, let $\gamma \in [0, 1)$, let TRUNCATE $\in \{\text{True}, \text{False}\}\ be\ a\ Boolean\ variable,\ and\ let\ (\mathcal{Y}_k)_{k=1:K} \subset S(\mathbb{R}^d)\ (recall\ Definition\ 3.1).$

The inputs of Algorithm 3 include a Boolean variable TRUNCATE \in {True, False} that indicates whether the probability measures ν_1, \ldots, ν_K in the outputs shall be truncated to $\mathcal{Y}_1, \ldots, \mathcal{Y}_K$. The following proposition presents the theoretical properties satisfied by ν_1, \ldots, ν_K both in the case where TRUNCATE = False and in the case where TRUNCATE = True. We refer the reader to the detailed discussion regarding when to use TRUNCATE = True and TRUNCATE = False after the proof of Proposition 5.2.

Proposition 5.2 (Synthetic generation of Wasserstein barycenter problem instance via Algorithm 3). Let the inputs of Algorithm 3 satisfy Setting 5.1, and let $(\nu_k, T_k)_{k=1:K}$ be the outputs of Algorithm 3. Moreover, let $(\overline{\lambda}_{\widetilde{k}})_{\widetilde{k}=1:\widetilde{K}}$ satisfy the condition in Line 5, let $(\beta_{-\widetilde{k}}, \beta_{\widetilde{k}})_{\widetilde{k}=1:\widetilde{K}}$ be defined in Line 8, and let the functions $\widetilde{\varphi}_1, \ldots, \widetilde{\varphi}_{\widetilde{K}}, \widetilde{\varphi}_{-1}, \ldots, \widetilde{\varphi}_{-\widetilde{K}}, \varphi_1, \ldots, \varphi_K, \varphi_1^*, \ldots, \varphi_K^* : \mathbb{R}^d \to \mathbb{R}$ be defined as follows:

$$\widetilde{\varphi}_{\widetilde{k}}(\boldsymbol{x}) := \theta_{\widetilde{k}} \log \left(\sum_{j=1}^{n_{\widetilde{k}}} \exp \left(\frac{1}{\theta_{\widetilde{k}}} \left(g_{\widetilde{k},j} + \langle \boldsymbol{y}_{\widetilde{k},j}, \boldsymbol{x} \rangle - \frac{1}{2} \| \boldsymbol{y}_{\widetilde{k},j} \|^2 \right) \right) \right) + \frac{\lambda_{\widetilde{k}}}{2} \| \boldsymbol{x} \|^2 \qquad \forall \boldsymbol{x} \in \mathbb{R}^d, \; \forall 1 \leq \widetilde{k} \leq \widetilde{K},$$

$$\widetilde{\varphi}_{-\widetilde{k}}(\boldsymbol{x}) := rac{\lambda_{\widetilde{k}}}{2} \|\boldsymbol{x}\|^2 - \widetilde{\varphi}_{\widetilde{k}}(\boldsymbol{x}) \qquad \qquad \forall \boldsymbol{x} \in \mathbb{R}^d, \ \forall 1 \leq \widetilde{k} \leq \widetilde{K},$$

$$\begin{split} \varphi_k(\boldsymbol{x}) &:= \left(\sum_{i \in \Phi^{-1}(k)} \beta_i \widetilde{\varphi}_i(\boldsymbol{x})\right) + \gamma \left\langle \frac{1}{2} \mathbf{A}_k \boldsymbol{x} + \boldsymbol{b}_k, \boldsymbol{x} \right\rangle & \forall \boldsymbol{x} \in \mathbb{R}^d, \ \forall 1 \le k \le K, \\ \varphi_k^*(\boldsymbol{y}) &:= \sup_{\boldsymbol{x} \in \mathbb{R}^d} \left\{ \langle \boldsymbol{y}, \boldsymbol{x} \rangle - \varphi_k(\boldsymbol{x}) \right\} & \forall \boldsymbol{y} \in \mathbb{R}^d, \ \forall 1 \le k \le K. \end{split}$$

Then, in the case where TRUNCATE = False, the following statement holds.

(i) $\nu_k \in \mathcal{P}_{2,ac}(\mathbb{R}^d)$ for k = 1, ..., K and $\bar{\mu}$ is the unique \mathcal{W}_2 -barycenter of $\nu_1, ..., \nu_K$ with weights $w_1, ..., w_K$.

In the case where TRUNCATE = True, let us assume in addition that $\bar{\mu} \in \mathcal{M}^q_{\text{full}}(\mathbb{R}^d)$ for $q \in \mathbb{N}_0$. Furthermore, for any $\mathcal{Y}_1, \ldots, \mathcal{Y}_K \in \mathcal{S}(\mathbb{R}^d)$, let $(\epsilon_k^{(1)}(\mathcal{Y}_k), \epsilon_k^{(2)}(\mathcal{Y}_k))_{k=1:K} \subset \mathbb{R}_+$ and $\epsilon(\mathcal{Y}_1, \ldots, \mathcal{Y}_K) \in \mathbb{R}_+$ be defined as follows:

$$\begin{aligned} \epsilon_k^{(1)}(\mathcal{Y}_k) &:= \int_{\mathbb{R}^d} 2\|\boldsymbol{x}\|^2 \Big(\frac{1-T_k \sharp \bar{\mu}(\mathcal{Y}_k)}{T_k \sharp \bar{\mu}(\mathcal{Y}_k)} + \mathbb{1}_{\mathcal{Y}_k^c}(\boldsymbol{x}) \Big) T_k \sharp \bar{\mu}(\mathrm{d}\boldsymbol{x}) & \forall 1 \le k \le K, \\ \epsilon_k^{(2)}(\mathcal{Y}_k) &:= \int_{\mathbb{R}^d} \left\| \|\boldsymbol{x}\|^2 - 2\varphi_k^*(\boldsymbol{x}) \right\| \Big(\frac{1-T_k \sharp \bar{\mu}(\mathcal{Y}_k)}{T_k \sharp \bar{\mu}(\mathcal{Y}_k)} + \mathbb{1}_{\mathcal{Y}_k^c}(\boldsymbol{x}) \Big) T_k \sharp \bar{\mu}(\mathrm{d}\boldsymbol{x}) & \forall 1 \le k \le K, \end{aligned}$$

$$\epsilon(\mathcal{Y}_1,\ldots,\mathcal{Y}_K) := \sum_{k=1}^{\infty} w_k \Big(2\mathcal{W}_2(\bar{\mu},T_k \sharp \bar{\mu}) \epsilon_k^{(1)}(\mathcal{Y}_k)^{\frac{1}{2}} + \epsilon_k^{(1)}(\mathcal{Y}_k) + \epsilon_k^{(2)}(\mathcal{Y}_k) \Big),$$

where $\mathcal{Y}_k^c := \mathbb{R}^d \setminus \mathcal{Y}_k$ for $k = 1, \dots, K$. Then, the following statements hold.

(ii) $\nu_k \in \mathcal{M}^q(\mathbb{R}^d)$ for $k = 1, \dots, K$ and

$$V(\bar{\mu}) \leq \inf_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \left\{ V(\mu) \right\} + \epsilon(\mathcal{Y}_1, \dots, \mathcal{Y}_K),$$
(5.1)

where $V(\cdot)$ is defined in (1.2) with respect to the weights w_1, \ldots, w_K in the inputs of Algorithm 3. In particular, ν_1, \ldots, ν_K satisfy Setting 3.6.

(iii) For k = 1, ..., K, let $(\mathcal{Y}_{k,r})_{r \in \mathbb{N}}$ be a family of increasing sets satisfying the conditions in Assumption 3.3. Then, it holds that $\lim_{r\to\infty} \epsilon(\mathcal{Y}_{1,r}, ..., \mathcal{Y}_{K,r}) = 0$.

Proof of Proposition 5.2. For $\tilde{k} = 1, ..., \tilde{K}$, it follows from Line 6, Lemma 4.2, and (4.4) in the proof of Lemma 4.2 that $\nabla \tilde{\varphi}_{\tilde{k}} = \tilde{T}_{\tilde{k}}$ and $\tilde{\varphi}_{\tilde{k}} \in \mathfrak{C}^{\infty}_{\underline{\lambda}_{\tilde{k}},\overline{\lambda}_{\tilde{k}}-\underline{\lambda}_{\tilde{k}}}(\mathbb{R}^d)$ where $\overline{\lambda}_{\tilde{k}} - \underline{\lambda}_{\tilde{k}} > \underline{\lambda}_{\tilde{k}}$. Moreover, it follows from Line 6 that $\nabla \tilde{\varphi}_{-\tilde{k}} = \tilde{T}_{-\tilde{k}}$. Observe that $\underline{\lambda}_{\tilde{k}} \mathbf{I}_d \preceq \nabla^2 \tilde{\varphi}_{-\tilde{k}}(\mathbf{x}) \preceq (\overline{\lambda}_{\tilde{k}} - \underline{\lambda}_{\tilde{k}}) \mathbf{I}_d$ for all $\mathbf{x} \in \mathbb{R}^d$, which implies that $\tilde{\varphi}_{-\tilde{k}} \in \mathfrak{C}^{\infty}_{\underline{\lambda}_{\tilde{k}},\overline{\lambda}_{\tilde{k}}-\underline{\lambda}_{\tilde{k}}}(\mathbb{R}^d)$. Subsequently, since $(\tilde{\varphi}_{\tilde{k}})_{\tilde{k}=1:\tilde{K}}$ and $(\tilde{\varphi}_{-\tilde{k}})_{\tilde{k}=1:\tilde{K}}$ are all infinitely differentiable, smooth, and strongly convex, and since $(\beta_{\tilde{k}})_{\tilde{k}=1:\tilde{K}} \subset (0,\infty), (\beta_{-\tilde{k}})_{\tilde{k}=1:\tilde{K}} \subset (0,\infty), \gamma \in [0,1), (\mathbf{A}_k)_{k=1:K} \subset \mathbb{S}^d_{++}$, it

holds by definition that, for k = 1, ..., K, $\varphi_k \in \mathfrak{C}^{\infty}_{\underline{\zeta}_k, \overline{\zeta}_k}(\mathbb{R}^d)$ for some $0 < \underline{\zeta}_k < \overline{\zeta}_k < \infty$. Hence, since Line 10 implies $T_k = \nabla \varphi_k$, it holds that T_k is $\overline{\zeta}_k$ -Lipschitz continuous and thus belongs to $C_{\text{lin}}(\mathbb{R}^d, \mathbb{R}^d)$. Furthermore, since Line 8 implies $w_{\Phi(\widetilde{k})}\beta_{\widetilde{k}} = w_{\Phi(-\widetilde{k})}\beta_{-\widetilde{k}}$ for $\widetilde{k} = 1, ..., \widetilde{K}$ and $\sum_{\widetilde{k}=1}^{\widetilde{K}} w_{\Phi(-\widetilde{k})}\beta_{-\widetilde{k}} = 1 - \gamma$, and since $\sum_{k=1}^{K} w_k \mathbf{A}_k = \mathbf{I}_d, \sum_{k=1}^{K} w_k \mathbf{b}_k = \mathbf{0}_d$ by assumption, we get

$$\sum_{k=1}^{K} w_{k} \varphi_{k}(\boldsymbol{x}) = \left(\sum_{k=1}^{K} w_{k} \sum_{i \in \Phi^{-1}(k)} \beta_{i} \widetilde{\varphi}_{i}(\boldsymbol{x})\right) + \gamma \sum_{k=1}^{K} \left\langle \frac{1}{2} w_{k} \mathbf{A}_{k} \boldsymbol{x} + w_{k} \boldsymbol{b}_{k}, \boldsymbol{x} \right\rangle$$

$$= \sum_{\widetilde{k}=1}^{\widetilde{K}} \left(w_{\Phi(\widetilde{k})} \beta_{\widetilde{k}} \widetilde{\varphi}_{\widetilde{k}}(\boldsymbol{x}) + w_{\Phi(-\widetilde{k})} \beta_{-\widetilde{k}} \widetilde{\varphi}_{-\widetilde{k}}(\boldsymbol{x}) \right) + \frac{\gamma}{2} \|\boldsymbol{x}\|^{2}$$

$$= \left(\sum_{\widetilde{k}=1}^{\widetilde{K}} \left(w_{\Phi(\widetilde{k})} \beta_{\widetilde{k}} - w_{\Phi(-\widetilde{k})} \beta_{-\widetilde{k}} \right) \widetilde{\varphi}_{\widetilde{k}}(\boldsymbol{x}) \right) + \frac{1}{2} \left(\sum_{\widetilde{k}=1}^{\widetilde{K}} w_{\Phi(-\widetilde{k})} \beta_{-\widetilde{k}} \overline{\lambda}_{\widetilde{k}} \right) \|\boldsymbol{x}\|^{2} + \frac{\gamma}{2} \|\boldsymbol{x}\|^{2}$$

$$= \frac{1-\gamma}{2} \|\boldsymbol{x}\|^{2} + \frac{\gamma}{2} \|\boldsymbol{x}\|^{2} = \frac{1}{2} \|\boldsymbol{x}\|^{2}.$$
(5.2)

Now, let us consider the case where TRUNCATE = False and prove statement (i). For k = 1, ..., K, since $T_k = \nabla \varphi_k \in \mathcal{C}_{\text{lin}}(\mathbb{R}^d, \mathbb{R}^d)$ for $\varphi_k \in \mathfrak{C}_{\underline{\zeta}_k, \overline{\zeta}_k}^{\infty}(\mathbb{R}^d)$, and since $\bar{\mu} \in \mathcal{P}_{2,\text{ac}}(\mathbb{R}^d)$ by assumption, it follows from Line 14 and Proposition 3.5(ii) that $\nu_k \in \mathcal{P}_{2,\text{ac}}(\mathbb{R}^d)$. Subsequently, the uniqueness of the \mathcal{W}_2 -barycenter of ν_1, \ldots, ν_K with weights w_1, \ldots, w_K follows directly from Theorem 2.3. In the following, we will prove that $\bar{\mu}$ is the unique \mathcal{W}_2 -barycenter of ν_1, \ldots, ν_K with weights w_1, \ldots, w_K through Brenier's theorem (Theorem 2.4); see also [13, Appendix C.2]. Notice that $[I_d, \nabla \varphi_k] \sharp \bar{\mu} \in \Pi(\bar{\mu}, \nu_k)$ is the unique optimal coupling between $\bar{\mu}$ and ν_k (where $I_d : \mathbb{R}^d \to \mathbb{R}^d$ denotes the identity map on \mathbb{R}^d) by Theorem 2.4, and it thus follows from (5.2) that

$$\sum_{k=1}^{K} w_{k} \mathcal{W}_{2}(\bar{\mu}, \nu_{k})^{2} = \sum_{k=1}^{K} w_{k} \left(\int_{\mathbb{R}^{d}} \|\boldsymbol{x}\|^{2} - 2\varphi_{k}(\boldsymbol{x}) \bar{\mu}(\mathrm{d}\boldsymbol{x}) + \int_{\mathbb{R}^{d}} \|\boldsymbol{y}\|^{2} - 2\varphi_{k}^{*}(\boldsymbol{y}) \nu_{k}(\mathrm{d}\boldsymbol{y}) \right)$$
$$= \int_{\mathbb{R}^{d}} \|\boldsymbol{x}\|^{2} - 2\left(\sum_{k=1}^{K} w_{k} \varphi_{k}(\boldsymbol{x})\right) \bar{\mu}(\mathrm{d}\boldsymbol{x}) + \left(\sum_{k=1}^{K} w_{k} \int_{\mathbb{R}^{d}} \|\boldsymbol{y}\|^{2} - 2\varphi_{k}^{*}(\boldsymbol{y}) \nu_{k}(\mathrm{d}\boldsymbol{y})\right)$$
$$= \sum_{k=1}^{K} w_{k} \int_{\mathbb{R}^{d}} \|\boldsymbol{y}\|^{2} - 2\varphi_{k}^{*}(\boldsymbol{y}) \nu_{k}(\mathrm{d}\boldsymbol{y}).$$
(5.3)

On the other hand, since $\varphi_k(x) + \varphi_k^*(y) \ge \langle x, y \rangle$ for all $x, y \in \mathbb{R}^d$, it holds that

$$\sum_{k=1}^{K} w_{k} \mathcal{W}_{2}(\mu, \nu_{k})^{2} \geq \sum_{k=1}^{K} w_{k} \left(\int_{\mathbb{R}^{d}} \|\boldsymbol{x}\|^{2} - 2\varphi_{k}(\boldsymbol{x}) \,\mu(\mathrm{d}\boldsymbol{x}) + \int_{\mathbb{R}^{d}} \|\boldsymbol{y}\|^{2} - 2\varphi_{k}^{*}(\boldsymbol{y}) \,\nu_{k}(\mathrm{d}\boldsymbol{y}) \right)$$

$$= \int_{\mathbb{R}^{d}} \|\boldsymbol{x}\|^{2} - 2\left(\sum_{k=1}^{K} w_{k} \varphi_{k}(\boldsymbol{x})\right) \,\mu(\mathrm{d}\boldsymbol{x}) + \left(\sum_{k=1}^{K} w_{k} \int_{\mathbb{R}^{d}} \|\boldsymbol{y}\|^{2} - 2\varphi_{k}^{*}(\boldsymbol{y}) \,\nu_{k}(\mathrm{d}\boldsymbol{y})\right)$$

$$= \sum_{k=1}^{K} w_{k} \int_{\mathbb{R}^{d}} \|\boldsymbol{y}\|^{2} - 2\varphi_{k}^{*}(\boldsymbol{y}) \,\nu_{k}(\mathrm{d}\boldsymbol{y}) \qquad \forall \mu \in \mathcal{P}_{2}(\mathbb{R}^{d}).$$
(5.4)

Combining (5.3) and (5.4) verifies that $\bar{\mu}$ is indeed the unique W_2 -barycenter of ν_1, \ldots, ν_K with weights w_1, \ldots, w_K . This completes the proof of statement (i).

Next, let us consider the case where TRUNCATE = True and prove statements (ii) and (iii). For k = 1, ..., K, since $T_k = \nabla \varphi_k \in \mathcal{C}_{\text{lin}}(\mathbb{R}^d, \mathbb{R}^d)$ for $\varphi_k \in \mathfrak{C}_{\underline{\zeta}_k, \overline{\zeta}_k}^{\infty}(\mathbb{R}^d)$, and since $\bar{\mu} \in \mathcal{M}_{\text{full}}^q(\mathbb{R}^d)$ by assumption, it follows from Proposition 3.5(iii) that $T_k \sharp \bar{\mu} \in \mathcal{M}_{\text{full}}^q(\mathbb{R}^d)$. Moreover, since $\mathcal{Y}_k \in \mathcal{S}(\mathbb{R}^d)$, it follows from Line 12 and Proposition 3.5(i) that $\nu_k := (T_k \sharp \bar{\mu})|_{\mathcal{Y}_k} \in \mathcal{M}^q(\mathbb{R}^d)$. To show that (5.1) holds, let us first derive an upper bound for $\mathcal{W}_2(T_k \sharp \bar{\mu}, \nu_k)^2$ for $k = 1, \ldots, K$. For every $k = 1, \ldots, K$, we define $\check{\nu}_k := (T_k \sharp \bar{\mu})|_{\mathcal{Y}_k^c}$, let $\pi_{k,1} := (T_k \sharp \bar{\mu})|_{\mathcal{Y}_k^c}$. $[I_d, I_d] \sharp \nu_k, \text{ let } \pi_{k,2} \in \Pi(\check{\nu}_k, \nu_k) \text{ be arbitrary, and let } \pi_k := T_k \sharp \bar{\mu}(\mathcal{Y}_k) \pi_{k,1} + (1 - T_k \sharp \bar{\mu}(\mathcal{Y}_k)) \pi_{k,2} \in \mathcal{P}(\mathbb{R}^d \times \mathbb{R}^d).$ One checks that $\pi_k \in \Pi(T_k \sharp \bar{\mu}, \nu_k)$. Consequently, we get

$$\begin{aligned} \mathcal{W}_{2}(T_{k}\sharp\bar{\mu},\nu_{k})^{2} &\leq \int_{\mathbb{R}^{d}\times\mathbb{R}^{d}} \|\boldsymbol{x}-\boldsymbol{y}\|^{2} \pi_{k}(\mathrm{d}\boldsymbol{x},\mathrm{d}\boldsymbol{y}) \\ &= T_{k}\sharp\bar{\mu}(\mathcal{Y}_{k}) \int_{\mathbb{R}^{d}} \|\boldsymbol{x}-\boldsymbol{x}\|^{2} \nu_{k}(\mathrm{d}\boldsymbol{x}) + (1-T_{k}\sharp\bar{\mu}(\mathcal{Y}_{k})) \int_{\mathbb{R}^{d}\times\mathbb{R}^{d}} \|\boldsymbol{x}-\boldsymbol{y}\|^{2} \pi_{k,2}(\mathrm{d}\boldsymbol{x},\mathrm{d}\boldsymbol{y}) \\ &\leq (1-T_{k}\sharp\bar{\mu}(\mathcal{Y}_{k})) \int_{\mathbb{R}^{d}} 2\|\boldsymbol{x}\|^{2} \check{\nu}_{k}(\mathrm{d}\boldsymbol{x}) + (1-T_{k}\sharp\bar{\mu}(\mathcal{Y}_{k})) \int_{\mathbb{R}^{d}} 2\|\boldsymbol{y}\|^{2} \nu_{k}(\mathrm{d}\boldsymbol{y}) \\ &= \int_{\mathbb{R}^{d}} 2\|\boldsymbol{x}\|^{2} \mathbb{1}_{\mathcal{Y}_{k}^{c}}(\boldsymbol{x}) T_{k}\sharp\bar{\mu}(\mathrm{d}\boldsymbol{x}) + \frac{1-T_{k}\sharp\bar{\mu}(\mathcal{Y}_{k})}{T_{k}\sharp\bar{\mu}(\mathcal{Y}_{k})} \int_{\mathbb{R}^{d}} 2\|\boldsymbol{y}\|^{2} \mathbb{1}_{\mathcal{Y}_{k}}(\boldsymbol{y}) T_{k}\sharp\bar{\mu}(\mathrm{d}\boldsymbol{y}) \\ &\leq \int_{\mathbb{R}^{d}} 2\|\boldsymbol{x}\|^{2} \Big(\frac{1-T_{k}\sharp\bar{\mu}(\mathcal{Y}_{k})}{T_{k}\sharp\bar{\mu}(\mathcal{Y}_{k})} + \mathbb{1}_{\mathcal{Y}_{k}^{c}}(\boldsymbol{x})\Big) T_{k}\sharp\bar{\mu}(\mathrm{d}\boldsymbol{x}) = \epsilon_{k}^{(1)}(\mathcal{Y}_{k}) \qquad \forall 1 \leq k \leq K. \end{aligned}$$

This leads to the following inequality:

$$V(\bar{\mu}) = \sum_{k=1}^{K} w_k \mathcal{W}_2(\bar{\mu}, \nu_k)^2$$

$$\leq \sum_{k=1}^{K} w_k \left(\mathcal{W}_2(\bar{\mu}, T_k \sharp \bar{\mu})^2 + 2\mathcal{W}_2(\bar{\mu}, T_k \sharp \bar{\mu}) \mathcal{W}_2(T_k \sharp \bar{\mu}, \nu_k) + \mathcal{W}_2(T_k \sharp \bar{\mu}, \nu_k)^2 \right)$$

$$\leq \left(\sum_{k=1}^{K} w_k \mathcal{W}_2(\bar{\mu}, T_k \sharp \bar{\mu})^2 \right) + \left(\sum_{k=1}^{K} w_k \left(2\mathcal{W}_2(\bar{\mu}, T_k \sharp \bar{\mu}) \epsilon_k^{(1)}(\mathcal{Y}_k)^{\frac{1}{2}} + \epsilon_k^{(1)}(\mathcal{Y}_k) \right) \right).$$
(5.5)

On the other hand, we have

$$\int_{\mathbb{R}^{d}} \left(\|\boldsymbol{y}\|^{2} - 2\varphi_{k}^{*}(\boldsymbol{y}) \right) T_{k} \sharp \bar{\mu}(\mathrm{d}\boldsymbol{y}) - \int_{\mathbb{R}^{d}} \left(\|\boldsymbol{y}\|^{2} - 2\varphi_{k}^{*}(\boldsymbol{y}) \right) \nu_{k}(\mathrm{d}\boldsymbol{y}) \\
\leq \int_{\mathbb{R}^{d}} \left\| \|\boldsymbol{y}\|^{2} - 2\varphi_{k}^{*}(\boldsymbol{y}) \right\| |T_{k} \sharp \bar{\mu} - \nu_{k}|(\mathrm{d}\boldsymbol{y}) \\
= \int_{\mathbb{R}^{d}} \left\| \|\boldsymbol{y}\|^{2} - 2\varphi_{k}^{*}(\boldsymbol{y}) \right\| \left| 1 - \frac{1}{T_{k} \sharp \bar{\mu}(\mathcal{Y}_{k})} \mathbb{1}_{\mathcal{Y}_{k}}(\boldsymbol{y}) \right| T_{k} \sharp \bar{\mu}(\mathrm{d}\boldsymbol{y}) \\
\leq \int_{\mathbb{R}^{d}} \left\| \|\boldsymbol{y}\|^{2} - 2\varphi_{k}^{*}(\boldsymbol{y}) \right\| \left(\frac{1 - T_{k} \sharp \bar{\mu}(\mathcal{Y}_{k})}{T_{k} \sharp \bar{\mu}(\mathcal{Y}_{k})} + \mathbb{1}_{\mathcal{Y}_{k}^{c}}(\boldsymbol{y}) \right) T_{k} \sharp \bar{\mu}(\mathrm{d}\boldsymbol{y}) = \epsilon_{k}^{(2)}(\mathcal{Y}_{k}) \quad \forall 1 \leq k \leq K.$$
(5.6)

Subsequently, combining (5.6) and (5.4) and then taking the infimum over all $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ leads to

$$\inf_{\boldsymbol{\mu}\in\mathcal{P}_{2}(\mathbb{R}^{d})}\left\{V(\boldsymbol{\mu})\right\} \geq \left(\sum_{k=1}^{K} w_{k} \int_{\mathbb{R}^{d}} \|\boldsymbol{y}\|^{2} - 2\varphi_{k}^{*}(\boldsymbol{y}) T_{k} \sharp \bar{\boldsymbol{\mu}}(\mathrm{d}\boldsymbol{y})\right) - \left(\sum_{k=1}^{K} w_{k} \epsilon_{k}^{(2)}(\boldsymbol{\mathcal{Y}}_{k})\right).$$
(5.7)

Furthermore, since $[I_d, \nabla \varphi_k] \sharp \bar{\mu} \in \Pi(\bar{\mu}, \nu_k)$ is the unique optimal coupling between $\bar{\mu}$ and $T_k \sharp \bar{\mu}$ by Brenier's theorem (Theorem 2.4), it holds that

$$\sum_{k=1}^{K} w_k \mathcal{W}_2(\bar{\mu}, T_k \sharp \bar{\mu})^2 = \sum_{k=1}^{K} w_k \left(\int_{\mathbb{R}^d} \|\boldsymbol{x}\|^2 - 2\varphi_k(\boldsymbol{x}) \,\bar{\mu}(\mathrm{d}\boldsymbol{x}) + \int_{\mathbb{R}^d} \|\boldsymbol{y}\|^2 - 2\varphi_k^*(\boldsymbol{y}) \,T_k \sharp \bar{\mu}(\mathrm{d}\boldsymbol{y}) \right)$$

$$= \sum_{k=1}^{K} w_k \int_{\mathbb{R}^d} \|\boldsymbol{y}\|^2 - 2\varphi_k^*(\boldsymbol{y}) \,T_k \sharp \bar{\mu}(\mathrm{d}\boldsymbol{y}).$$
 (5.8)

Now, combining (5.5), (5.7), and (5.8) proves (5.1) and completes the proof of statement (ii). Lastly, for k = 1, ..., K, observe that φ_k^* is $\underline{\zeta}_k^{-1}$ -smooth and $\overline{\zeta}_k^{-1}$ -strongly convex by the duality be-tween smooth convex functions and strongly convex functions (see, e.g., the equivalence between (a) and (e) in [56, Proposition 12.60]). Hence, φ_k^* is bounded from below by some constant and dominated from above by some quadratic function, e.g., by $\mathbb{R}^d \ni \mathbf{x} \mapsto \frac{1}{\underline{\zeta}_k} \|\mathbf{x}\|^2 + C \in \mathbb{R}$ for sufficiently large C > 0. Consequently, for k = 1, ..., K, the property $\bigcup_{r \in \mathbb{N}} \mathcal{Y}_{k,r} = \mathbb{R}^d$ and Lebesgue's dominated convergence theorem imply that $\lim_{r\to\infty} \epsilon_k^{(1)}(\mathcal{Y}_{k,r}) = \lim_{r\to\infty} \epsilon_k^{(2)}(\mathcal{Y}_{k,r}) = 0$ for k = 1, ..., K, and we thus get $\lim_{r\to\infty} \epsilon(\mathcal{Y}_{1,r}, ..., \mathcal{Y}_{K,r}) = 0$, which proves statement (iii). The proof is now complete. \Box

We provide here the motivation as well as a summary of the results in Proposition 5.2. Proposition 5.2(i) shows that, for any $\bar{\mu} \in \mathcal{P}_{2,ac}(\mathbb{R}^d)$, Algorithm 3 with TRUNCATE set to False constructs $\nu_1, \ldots, \nu_K \in$ $\mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$ such that $\bar{\mu}$ is their unique \mathcal{W}_2 -barycenter with the user-specified weights w_1,\ldots,w_K . However, ν_1, \ldots, ν_K constructed this way do not necessarily satisfy Setting 3.6 for Algorithm 2. In fact, the ν_1, \ldots, ν_K may not belong to $\mathcal{M}^q(\mathbb{R}^d)$ even if $\bar{\mu} \in \mathcal{M}^q(\mathbb{R}^d)$ is imposed, since the supports of $T_1 \sharp \bar{\mu}, \ldots, T_K \sharp \bar{\mu}$ are not necessarily convex. Proposition 5.2(ii) on the other hand, shows that for any $\bar{\mu} \in \mathcal{M}^q_{\text{full}}(\mathbb{R}^d)$, Algorithm 3 with TRUNCATE set to True constructs ν_1, \ldots, ν_K satisfying Setting 3.6, and $\bar{\mu}$ approximates their unique \mathcal{W}_2 -barycenter with the user-specified weights w_1, \ldots, w_K , in the sense that $\bar{\mu}$ approximately solves the minimization problem $\inf_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \{V(\mu)\}$ that characterizes the \mathcal{W}_2 -barycenter, where the suboptimality of $\bar{\mu}$ is bounded by $\epsilon(\mathcal{Y}_1, \ldots, \mathcal{Y}_K) > 0$. Furthermore, Proposition 5.2(iii) shows that one may control the sub-optimality of $\bar{\mu}$ by letting $\mathcal{Y}_k \leftarrow \mathcal{Y}_{k,r}$ for $k = 1, \ldots, K$ in the inputs of Algorithm 3, where $(\mathcal{Y}_{1,r})_{r\in\mathbb{N}},\ldots,(\mathcal{Y}_{K,r})_{r\in\mathbb{N}}$ are families of increasing sets satisfying the conditions in Assumption 3.3. Subsequently, the sub-optimality of $\bar{\mu}$ can be controlled to be arbitrarily small by choosing sufficiently large $r \in \mathbb{N}$. Hence, when $r \in \mathbb{N}$ is large, $\overline{\mu}$ is a highly accurate approximate \mathcal{W}_2 -barycenter of the generated measures ν_1, \ldots, ν_K with user-specified weights w_1, \ldots, w_K . As such, $V(\bar{\mu})$ can be treated as an approximate lower bound when we quantitatively analyze the empirical approximation error of any W_2 -barycenter algorithm using the generated measures ν_1, \ldots, ν_K .

Remark 5.3. Under the settings of Proposition 5.2(iii) where $(\mathcal{Y}_{1,r})_{r\in\mathbb{N}}, \ldots, (\mathcal{Y}_{K,r})_{r\in\mathbb{N}}$ are families of increasing sets satisfying the conditions in Assumption 3.3, the result of Le Gouic and Loubes [39, Proposition 6] about the stability of Wasserstein barycenter can be used to conclude that the unique \mathcal{W}_2 -barycenter of ν_1, \ldots, ν_K generated by Algorithm 3 with TRUNCATE \leftarrow True and $\mathcal{Y}_k \leftarrow \mathcal{Y}_{k,r}$ for $k = 1, \ldots, K$ converges in \mathcal{W}_2 to $\bar{\mu}$ as $r \to \infty$. However, we are unable to get any quantitative bound on their \mathcal{W}_2 -distance due to the non-compactness of the supports of $T_1 \sharp \bar{\mu}, \ldots, T_K \sharp \bar{\mu}$.

On the practical side, assuming that independent samples from $\bar{\mu} \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$ in the inputs of Algorithm 3 can be efficiently generated, one can efficiently generate independent samples from ν_1, \ldots, ν_K in the outputs of Algorithm 3 as follows. First, one generates $S \in \mathbb{N}$ independent samples $Z^{[1]}, \ldots, Z^{[S]}$ from $\bar{\mu}$ and then computes their images $T_k(Z^{[1]}), \ldots, T_k(Z^{[S]})$ under T_k for each $k = 1, \ldots, K$. Subsequently, if TRUNCATE is set to False, then $\{T_k(Z^{[1]}), \ldots, T_k(Z^{[S]})\}$ are S independent samples from ν_k , for $k = 1, \ldots, K$. If TRUNCATE is set to True, then one performs an extra rejection step, where, for $k = 1, \ldots, K$, the samples in $\{T_k(Z^{[1]}), \ldots, T_k(Z^{[S]})\}$ that do not belong to \mathcal{Y}_k are rejected and discarded. In this way, $\{T_k(Z^{[1]}), \ldots, T_k(Z^{[S]})\} \cap \mathcal{Y}_k$ are independent samples from ν_k , for $k = 1, \ldots, K$. When one chooses the sets $\mathcal{Y}_1, \ldots, \mathcal{Y}_K$ in the inputs of Algorithm 3 to be sufficiently "large", e.g., by letting $\mathcal{Y}_1 = \cdots = \mathcal{Y}_K = \bar{B}(\mathbf{0}_d, r)$ for sufficiently large r, the rejection ratio will be close to 0 and the impact of the truncation in Line 12 of Algorithm 3 will be unnoticeable in practice.

6. NUMERICAL EXPERIMENTS

In this section, we apply Algorithm 3 in Section 5 for generating non-trivial problem instances to assess Wasserstein barycenter algorithms, and present our numerical experiments to highlight the efficacy of our proposed stochastic fixed-point algorithm (Algorithm 2) when deployed with the modified entropic OT map estimators in Proposition 4.1. Specifically, under identical problem instances, we compare our results with two state-of-the-art free-support methods for approximating W_2 -barycenters proposed by Korotin, Egiazarian, Li, and Burnaev [38] and Fan, Taghvaei, and Chen [26] respectively, which employ neural networks (NNs) to parametrize Brenier potentials and optimal transportation maps while employing generative neural networks (GNNs) to parametrize the W_2 -barycenter. In our experiments, problem instances with d = 2 and d = 10 are considered.⁶ The Python implementation for our proposed algorithm and all codes for our numerical experiments can be accessed at https://github.com/CHENZeyi1101/WB_Algo.

⁶We remark that the task of finding the 2-dimensional Wasserstein barycenter has been widely witnessed in the area of computer vision and graphics, while high-dimensional barycenter approximation may be of particular interest in Bayesian inference, ensemble learning, data-driven decision-making, and many other applications.



FIGURE 6.1. Density visualizations of the approximate ground-truth W_2 -barycenter $\bar{\mu}$ and input measures ν_1, \ldots, ν_5 generated by Algorithm 3 in the 2-dimensional instance of our experiment.

Experimental settings. We present the inputs of Algorithm 3 in our experimental settings for generating synthetic problem instances following the notions in Setting 5.1. To begin, we pick an arbitrary $\bar{\mu} \in \mathcal{M}_{full}^q(\mathbb{R}^d)$ as the ground-truth barycenter measure, and fix $K \in \mathbb{N} \cap [2, \infty)$, $\tilde{K} \in \mathbb{N} \cap [2, \infty)$ with $2\tilde{K} \geq K$, as well as $\underline{\lambda}_{\tilde{k}} > 0$, $\theta_{\tilde{k}} > 0$, $n_{\tilde{k}} \in \mathbb{N}$ for $\tilde{k} = 1, \ldots, \tilde{K}$. Throughout the experiments, we consider barycenter problems with uniform weights; i.e., $w_1 = \ldots w_K = \frac{1}{K}$. In order to generate ν_1, \ldots, ν_K with non-trivial structures in Algorithm 3, we consider an arbitrary set of \tilde{K} auxiliary measures denoted by $\varkappa_1, \ldots, \varkappa_{\tilde{K}} \in \mathcal{P}_{2,\mathrm{ac}}(\mathbb{R}^d)$. We subsequently perform entropic OT map estimation between $\bar{\mu}$ and $\varkappa_{\tilde{k}}$ for $\tilde{k} = 1, \ldots, \tilde{K}$ to generate $(g_{\tilde{k},j}, y_{\tilde{k},j})_{j=1:n_{\tilde{k}}}$. To do so, we randomly generate $n_{\tilde{k}}$ independent samples $(x_{\tilde{k},j})_{j=1:n_{\tilde{k}}}$ from $\bar{\mu}$, randomly generate $n_{\tilde{k}}$ independent samples $(y_{\tilde{k},j})_{j=1:n_{\tilde{k}}}$, for $g(g_j^{0})_{j=1:n_{\tilde{k}}}$ from $\bar{\mu}$, randomly generate $n_{\tilde{k}}$ independent samples $(y_{\tilde{k},j})_{j=1:n_{\tilde{k}}}$ from $\omega(q_{\tilde{k},j})_{j=1:n_{\tilde{k}}}$ from $\omega(q_{\tilde{k},j})_{j=1:n_{\tilde{k}}}$ ($Y_j)_{j=1:n_{\tilde{k}}$ from $\varkappa_{\tilde{k}}$, and solve (4.1) with $m \leftarrow n_{\tilde{k}}$, $n \leftarrow n_{\tilde{k}}$, $(X_i)_{i=1:m} \leftarrow (x_{\tilde{k},i})_{i=1:n_{\tilde{k}}}$, $(Y_j)_{j=1:n_{\tilde{k}}} \in (y_{\tilde{k},j})_{j=1:n_{\tilde{k}}} \leftarrow (g_j^0)_{j=1:n}$. In this way, when $(\lambda_{\tilde{k}})_{\tilde{k}=1:\tilde{K}}$ are small, each $\tilde{T}_{\tilde{k}}$ defined in Line 6 of Algorithm 3 approximates the OT map $T_{\varkappa_{\tilde{k}}}^{\mu}$ from μ to $\varkappa_{\tilde{k}}$, and we are thus able to grant $(\tilde{T}_{\tilde{k}})_{\tilde{k}=1:\tilde{K}}$ and $(\tilde{T}_{-\tilde{k}})_{\tilde{k}=1:\tilde{K}}$ defined in Line 10 of Algorithm 3 acquire these non-affine structures from $(\tilde{T}_{\tilde{k}})_{\tilde{k}=1:\tilde{K}}$ and $(\tilde{T}_{-\tilde{k}})_{\tilde{k}=1:\tilde{K}}$. Moreover, we set TRUNCATE = True, and choose $(\mathcal{Y}_k)_{k=1:K}$ to be closed balls that are centered at the origin with sufficiently large radi

In the 2-dimensional case, we examine the performance of our algorithm in approximating the W_2 barycenter of K = 5 probability measures in $\mathcal{M}^q(\mathbb{R}^2)$. When generating our synthetic problem instance via Algorithm 3, we let the approximate ground-truth W_2 -barycenter $\bar{\mu}$ be a Gaussian mixture (GM) with 5 components. To generate the input measures ν_1, \ldots, ν_5 via Algorithm 3, we set $n_{\tilde{k}} = 1000$ and employ $\tilde{K} = 5$ GMs, namely $\varkappa_1, \ldots, \varkappa_5$, as the auxiliary measures. We show in Figure 6.1 the probability density function of $\bar{\mu}$ and the approximate densities of the input measures ν_1, \ldots, ν_5 via kernel density estimation (KDE) respectively from 2000 independent samples, where the color bars show the density scale. In the 10-dimensional case, we similarly set the approximate ground-truth measure and the auxiliary measures to be GMs each with 5 components, and set $K = \tilde{K} = 10$. For both cases, we start Algorithm 2 by setting the initial measure $\rho_0 = \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ to be the standard multivariate Gaussian distribution. Regarding other specific parameters (e.g., the sample sizes, the truncation radii, the regularization parameters, etc.) selected in each iteration of Algorithm 2, readers are referred to the aforementioned GitHub repository for further details.

Result analysis. Given a problem instance with a known approximate W_2 -barycenter generated by Algorithm 3, we are interested in two metrics to quantify the performance of each approximation algorithm, i.e., the accuracy of an approximately computed $\hat{\mu} \in \mathcal{P}_2(\mathbb{R}^d)$ relative to the (approximate) ground-truth barycenter $\bar{\mu}$: their difference in *V*-value, i.e., $|V(\hat{\mu}) - V(\bar{\mu})|$ (see Definition 1.1) with $w_1 = \cdots = w_K = \frac{1}{K}$, and their Wasserstein distance $W_2(\hat{\mu}, \bar{\mu})$. Since $\hat{\mu}$ and $\bar{\mu}$ are general continuous measures, the exact computation of the W_2 -distance between them is intractable (see, e.g., [62]). Therefore, we use the empirical W_2 -distance denoted by $\widehat{W}_{2,n}(\cdot, \cdot)$ as a proxy which is obtained by evaluating the associated *n*-sample empirical measures. In analogy, we evaluate the *V*-value using its discrete counterpart defined by $\widehat{V}_n(\cdot) := \frac{1}{K} \sum_{k=1}^K \widehat{W}_{2,n}(\cdot, \nu_k)^2$. In our experiments, $\hat{\mu}_t$ denotes the approximated W_2 -barycenter of ν_1, \ldots, ν_K at iteration *t* of our algorithm and the algorithm of Korotin et al. [38], or the approximated W_2 -barycenter at epoch *t* of the algorithm of Fan et al.



FIGURE 6.2. The evolution of $\widehat{V}_n(\widehat{\mu}_t)$ and $\widehat{W}_{2,n}(\widehat{\mu}_t, \overline{\mu})$ given by our proposed Algorithm 2 in the 2-dimensional and 10-dimensional problem instances.

[26], in order to respect their original notions. For each evaluation at $t \in \mathbb{N}_0$, we take n = 5000 independent samples respectively from $\hat{\mu}_t$ and $\bar{\mu}$ to evaluate $\hat{V}_n(\hat{\mu}_t) \approx V(\hat{\mu}_t)$ and $\widehat{W}_{2,n}(\hat{\mu}_t, \bar{\mu}) \approx \mathcal{W}_2(\hat{\mu}_t, \bar{\mu})$.

For both the 2-dimensional instance and the 10-dimensional instance, we compute 20 independent empirical approximations of $V(\bar{\mu})$ denoted by $(\hat{V}_n^{(i)}(\bar{\mu}))_{i=1:20}$, and evaluate their 10% trimmed mean and interquartile range (IQR). For each iteration or epoch t, we also compute 20 independent empirical approximations of $V(\hat{\mu}_t)$ denoted by $(\hat{V}_n^{(i)}(\hat{\mu}_t))_{i=1:20}$, and compute their 10% trimmed mean and IQR accordingly. Similarly, we examine the 10% trimmed mean and the IQR of 20 independent empirical approximations of $W_2(\hat{\mu}_t, \bar{\mu})$ denoted by $(\widehat{W}_{2,n}^{(i)}(\hat{\mu}_t, \bar{\mu}))_{i=1:20}$ for each iteration or epoch t. For all three examined algorithms, the numeric evolutions of $\widehat{V}_n(\hat{\mu}_t)$ and $\widehat{W}_{2,n}(\hat{\mu}_t, \bar{\mu})$ across iterations/epochs are evaluated in both 2-dimensional and 10-dimensional instances, and the results from the three examined algorithms are presented in Figures 6.2, 6.3 and 6.4 respectively. We emphasize that the algorithm-specific subplots therein cannot be combined and compared in a single figure due to incompatible iteration/epoch indices, and we follow the scales of iterations and epochs for the two benchmark algorithms as defined in their original settings.⁷

From Figure 6.2, one can empirically observe that our algorithm simultaneously reduces $\hat{V}_n(\hat{\mu}_t)$ and $\widehat{W}_{2,n}(\hat{\mu}_t, \bar{\mu})$ to nearly optimal levels within a few iterations, after which the approximated barycenter $\hat{\mu}_t$ remains close to the ground-truth barycenter $\bar{\mu}$. Between the two benchmark algorithms, Figure 6.3 implies that the algorithm of Korotin et al. [38] exhibits relatively poor performances, as the approximated $\hat{V}_n(\hat{\mu}_t)$ and $\widehat{W}_{2,n}(\hat{\mu}_t, \bar{\mu})$ therein tend to deviate from their optimal baselines in the 2-dimensional case, and fail to be optimized in the 10-dimensional case. This is possibly attributed to a misspecification of hyperparameters in their architecture, which may result in inferior parametrizations by the neural networks. The algorithm of Fan et al. [26], on the other hand, achieves a rapid descent in $\widehat{V}_n(\hat{\mu}_t)$ to near-optimality as witnessed in Figure 6.4, though it consumes more epochs for $\widehat{W}_{2,n}(\hat{\mu}_t, \bar{\mu})$ to reach near-optimality. A potential explanation for this phenomenon is that it can be more challenging for the generative neural network to learn the underlying geometric structure of the W_2 -barycenter despite its efficiency in optimizing the V-value. Moreover, it can be empirically observed in both dimensions that our Algorithm 2 attains smaller errors in $\widehat{V}_n(\widehat{\mu}_t)$ and $\widehat{W}_{2,n}(\widehat{\mu}_t, \bar{\mu})$ when compared with the algorithm of Fan et al. [26]. To gain further insights about the superior performance of our algorithm in the 2-dimensional problem instance, we sample from the measures $\widehat{\mu}_0$ and $\widehat{\mu}_1$ in Algorithm 2 which are visualized

⁷See https://github.com/iamalexkorotin/WassersteinIterativeNetworks for the open-source codes of Korotin et al. [38]. See https://github.com/sbyebss/Scalable-Wasserstein-Barycenter for the open-source codes of Fan et al. [26].



FIGURE 6.3. The evolution of $\widehat{V}_n(\widehat{\mu}_t)$ and $\widehat{W}_{2,n}(\widehat{\mu}_t, \overline{\mu})$ given by [38] in the 2-dimensional and 10-dimensional problem instances.



FIGURE 6.4. The evolution of $\widehat{V}_n(\widehat{\mu}_t)$ and $\widehat{W}_{2,n}(\widehat{\mu}_t, \overline{\mu})$ given by [26] in the 2-dimensional and 10-dimensional problem instances.

via KDE in Figure 6.5. The KDE heatmaps generated by 2000 independent samples from $\hat{\mu}_1$ and $\bar{\mu}$ then exhibit similar supports and density functions, which implies the accuracy of our algorithm in approximating the W_2 -barycenter simply within a single iteration.

Conclusion. We provide a stochastic fixed-point algorithm for approximately computing the W_2 -barycenter of continuous non-parametric measures along with its theoretical convergence guarantee. Compared to the state-of-the-art neural network based algorithms, our algorithm possesses remarkable advantages in terms of both accuracy and computational efficiency which are detailed as follows. Firstly, when estimating the OT maps, our algorithm is driven by the Sinkhorn's algorithm which is known to be highly efficient, while training neural networks amounts to propagating through involved architectures. Secondly, our method circumvents the need to tune sophisticated hyperparameters and avoids the model over-parametrization issues that are potentially present in neural networks, thus providing ease in parametrizing the underlying Brenier potentials and OT maps. Thirdly, empirical observations from our numerical experiments have revealed that our algorithm requires only a handful of iterations to achieve near-optimality. Lastly, our algorithm can be executed on standard CPUs without



FIGURE 6.5. Left - Visualizations of $\hat{\mu}_0$ and $\hat{\mu}_1$ from our proposed Algorithm 2, where 2000 independent samples from $\hat{\mu}_0$ and $\hat{\mu}_1$ are plotted on the planes Y = 0 and Y = 50 respectively, and the KDE heatmap generated from $\hat{\mu}_1$ is plotted on Y = 100. Right - The KDE heatmap generated by 2000 independent samples from $\bar{\mu}$.

the need for GPUs or high-performance computing hardware, while also permitting hardware acceleration (see, e.g., [27] and the GeomLoss library therein).

ACKNOWLEDGEMENTS

ZC gratefully acknowledges the support from Nanyang Technological University under the URECA Undergraduate Research Programme, and the financial support from the INSEAD PhD Fellowship. AN and QX gratefully acknowledge the financial support by the MOE AcRF Tier 2 Grant *MOE-T2EP20222-0013*.

REFERENCES

- [1] M. Agueh and G. Carlier. Barycenters in the Wasserstein space. *SIAM Journal on Mathematical Analysis*, 43(2):904–924, 2011.
- [2] J. M. Altschuler and E. Boix-Adserà. Wasserstein barycenters are NP-hard to compute. SIAM Journal on Mathematics of Data Science, 4(1):179–203, 2022.
- [3] P. C. Álvarez-Esteban, E. Del Barrio, J. Cuesta-Albertos, and C. Matrán. A fixed-point approach to barycenters in Wasserstein space. *Journal of Mathematical Analysis and Applications*, 441(2):744–762, 2016.
- [4] L. Ambrosio, N. Gigli, and G. Savaré. *Gradient flows: in metric spaces and in the space of probability measures.* Springer Science & Business Media, 2008.
- [5] B. Amos, L. Xu, and J. Z. Kolter. Input convex neural networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 146–155. PMLR, 2017.
- [6] R. Bardenet, A. Doucet, and C. Holmes. On Markov chain Monte Carlo methods for tall data. *Journal of Machine Learning Research*, 18(47):1–43, 2017.
- [7] H. Bastani, D. Simchi-Levi, and R. Zhu. Meta dynamic pricing: Transfer learning across experiments. *Management Science*, 68(3):1865–1881, 2022.
- [8] G. Blekherman, P. A. Parrilo, and R. R. Thomas. *Semidefinite optimization and convex algebraic geometry*. Society for Industrial and Applied Mathematics, USA, 2012.
- [9] L. A. Caffarelli. A localization property of viscosity solutions to the Monge-Ampère equation and their strict convexity. *Annals of Mathematics*, 131(1):129–134, 1990.
- [10] L. A. Caffarelli. Some regularity properties of solutions of Monge-Ampère equation. Communications on Pure and Applied Mathematics, 44(8-9):965–969, 1991.
- [11] L. A. Caffarelli. The regularity of mappings with a convex potential. *Journal of the American Mathematical Society*, 5(1):99–104, 1992.
- [12] L. A. Caffarelli. Boundary regularity of maps with convex potentials–II. Annals of Mathematics, 144(3): 453–496, 1996.

- [13] S. Chewi, T. Maunu, P. Rigollet, and A. Stromme. Gradient descent algorithms for Bures-Wasserstein barycenters. In *Proceedings of 33rd Conference on Learning Theory*, volume 125, pages 1276–1304. PMLR, 2020.
- [14] S. Chewi, J. Niles-Weed, and P. Rigollet. Statistical optimal transport. Preprint arXiv:2407.18163, 2024.
- [15] L. Chizat, P. Roussillon, F. Léger, F.-X. Vialard, and G. Peyré. Faster Wasserstein distance estimation with the Sinkhorn divergence. In *Advances in Neural Information Processing Systems*, volume 33, pages 2257–2269, 2020.
- [16] S. Claici, E. Chien, and J. Solomon. Stochastic Wasserstein barycenters. In Proceedings of the 35th International Conference on Machine Learning, volume 80, pages 999–1008. PMLR, 2018.
- [17] S. Cohen, M. Arbel, and M. P. Deisenroth. Estimating barycenters of measures in high dimensions. *Preprint arXiv:2007.07105*, 2020.
- [18] A. D. D. Craik. Prehistory of Faà di Bruno's formula. *The American Mathematical Monthly*, 112(2): 119–130, 2005.
- [19] M. Curmei and G. Hall. Shape-constrained regression using sum of squares polynomials. *Operations Research*, 73(1):543–559, 2023.
- [20] M. Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In Advances in Neural Information Processing Systems, volume 26, 2013.
- [21] M. Cuturi, L. Meng-Papaxanthos, Y. Tian, C. Bunne, G. Davis, and O. Teboul. Optimal transport tools (OTT): A JAX toolbox for all things Wasserstein. *Preprint arXiv:2201.12324*, 2022.
- [22] N. Deb, P. Ghosal, and B. Sen. Rates of estimation of optimal transport maps using plug-in estimators via barycentric projections. In *Advances in Neural Information Processing Systems*, volume 34, pages 29736–29753, 2021.
- [23] P. Dognin, I. Melnyk, Y. Mroueh, J. Ross, C. D. Santos, and T. Sercu. Wasserstein barycenter model ensembling. *Preprint arXiv*:1902.04999, 2019.
- [24] A. L. Dontchev and R. T. Rockafellar. *Implicit functions and solution mappings: A view from variational analysis.* Springer Monographs in Mathematics. Springer, Dordrecht, 2009.
- [25] L. C. Evans. Partial Differential Equations, volume 19 of Graduate Studies in Mathematics. American Mathematical Society, 2nd edition, 2010.
- [26] J. Fan, A. Taghvaei, and Y. Chen. Scalable computations of Wasserstein barycenter via input convex neural networks. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pages 1571–1581. PMLR, 2021.
- [27] J. Feydy, P. Roussillon, A. Trouvé, and P. Gori. Fast and scalable optimal transport for brain tractograms. In *Medical Image Computing and Computer Assisted Intervention – MICCAI 2019*, volume 11767 of *Lecture Notes in Computer Science*, pages 636–644. Springer, 2019.
- [28] M. Fréchet. Les éléments aléatoires de nature quelconque dans un espace distancié. *Annales de l'institut Henri Poincaré*, 10:215–310, 1948.
- [29] A. Genevay, M. Cuturi, G. Peyré, and F. Bach. Stochastic optimization for large-scale optimal transport. In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- [30] N. Gigli. On Hölder continuity-in-time of the optimal transport map towards measures along a curve. *Proceedings of the Edinburgh Mathematical Society*, 54(2):401–409, 2011.
- [31] A. González-Sanz, L. De Lara, L. Béthune, and J.-M. Loubes. GAN estimation of Lipschitz optimal transport maps. *Preprint arXiv:2202.07965*, 2022.
- [32] F. F. Gunsilius. On the convergence rate of potentials of Brenier maps. *Econometric Theory*, 38(2): 381–417, 2022.
- [33] V. Gupta and N. Kallus. Data pooling in stochastic optimization. *Management Science*, 68(3):1595–1615, 2022.
- [34] J.-C. Hütter and P. Rigollet. Minimax estimation of smooth optimal transport maps. *The Annals of Statistics*, 49(2):1166–1194, 2021.
- [35] L. V. Kantorovich. On a problem of Monge. CR (Doklady) Acad. Sci. URSS (NS), 3:225-226, 1948.
- [36] H. Karcher. Riemannian center of mass and mollifier smoothing. *Communications on Pure and Applied Mathematics*, 30(5):509–541, 1977.
- [37] A. Korotin, L. Li, J. Solomon, and E. Burnaev. Continuous Wasserstein-2 barycenter estimation without minimax optimization. In *The 9th International Conference on Learning Representations*, 2021.

- [38] A. Korotin, V. Egiazarian, L. Li, and E. Burnaev. Wasserstein iterative networks for barycenter estimation. In *Advances in Neural Information Processing Systems*, volume 35, pages 15672–15686, 2022.
- [39] T. Le Gouic and J.-M. Loubes. Existence and consistency of Wasserstein barycenters. *Probab. Theory Related Fields*, 168(3-4):901–917, 2017.
- [40] L. Li, A. Genevay, M. Yurochkin, and J. M. Solomon. Continuous regularized Wasserstein barycenters. In Advances in Neural Information Processing Systems, volume 33, pages 17755–17765, 2020.
- [41] G. Luise, S. Salzo, M. Pontil, and C. Ciliberto. Sinkhorn barycenters with free support via Frank-Wolfe algorithm. In Advances in Neural Information Processing Systems, volume 32, 2019.
- [42] T. Manole, S. Balakrishnan, J. Niles-Weed, and L. Wasserman. Plugin estimation of smooth optimal transport maps. *The Annals of Statistics*, 52(3):966–998, 2024.
- [43] A. J. McNeil, R. Frey, and P. Embrechts. *Quantitative risk management: concepts, techniques and tools-revised edition*. Princeton University Press, 2015.
- [44] S. Minsker, S. Srivastava, L. Lin, and D. B. Dunson. Robust and scalable Bayes via a median of subset posterior measures. *The Journal of Machine Learning Research*, 18(1):4488–4527, 2017.
- [45] B. Muzellec and M. Cuturi. Generalizing point embeddings using the Wasserstein space of elliptical distributions. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- [46] B. Muzellec, A. Vacher, F. Bach, F.-X. Vialard, and A. Rudi. Near-optimal estimation of smooth transport maps with kernel sums-of-squares. *Preprint arXiv:2112.01907*, 2021.
- [47] Y. Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2004.
- [48] A. Neufeld and Q. Xiang. Feasible approximation of matching equilibria for large-scale matching for teams problems. *Preprint arXiv:2308.03550*, 2024.
- [49] F.-P. Paty, A. d'Aspremont, and M. Cuturi. Regularity as regularization: Smooth and strongly convex Brenier potentials in optimal transport. In *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics*, volume 108, pages 1222–1232. PMLR, 2020.
- [50] E. V. Petracou, A. Xepapadeas, and A. N. Yannacopoulos. Decision making under model uncertainty: Fréchet–Wasserstein mean preferences. *Management Science*, 68(2):1195–1211, 2022.
- [51] G. Peyré and M. Cuturi. Computational optimal transport: With applications to data science. *Foundations and Trends*® *in Machine Learning*, 11(5-6):355–607, 2019.
- [52] A.-A. Pooladian and J. Niles-Weed. Entropic estimation of optimal transport maps. *Preprint* arXiv:2109.12004, 2024.
- [53] J. Rabin, G. Peyré, J. Delon, and M. Bernot. Wasserstein barycenter and its application to texture mixing. In Scale Space and Variational Methods in Computer Vision: Third International Conference, SSVM 2011, Ein-Gedi, Israel, May 29–June 2, 2011, Revised Selected Papers 3, pages 435–446. Springer, 2012.
- [54] R. Ranjan and T. Gneiting. Combining probability forecasts. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 72(1):71–91, 2010.
- [55] R. T. Rockafellar. Convex Analysis: (PMS-28). Princeton university press, 1970.
- [56] R. T. Rockafellar and R. J.-B. Wets. Variational analysis, volume 317 of Grundlehren der mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]. Springer-Verlag, Berlin, 1998.
- [57] Y. Rychener, A. Esteban-Pérez, J. M. Morales, and D. Kuhn. Wasserstein distributionally robust optimization with heterogeneous data sources. *Preprint arXiv:2407.13582*, 2024.
- [58] F. Santambrogio. *Optimal transport for applied mathematicians*, volume 87 of *Progress in Nonlinear Differential Equations and their Applications*. Birkhäuser/Springer, Cham, 2015.
- [59] A. Spelta and N. Pecora. Wasserstein barycenter for link prediction in temporal networks. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 187(1):180–208, 2024.
- [60] S. Srivastava, V. Cevher, Q. Dinh, and D. Dunson. WASP: Scalable Bayes via barycenters of subset posteriors. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, volume 38, pages 912–920. PMLR, 2015.
- [61] S. Srivastava, C. Li, and D. B. Dunson. Scalable Bayes via barycenter in Wasserstein space. *Journal of Machine Learning Research*, 19(8):1–35, 2018.
- [62] B. Taşkesen, S. Shafieezadeh-Abadeh, and D. Kuhn. Semi-discrete optimal transport: hardness, regularization and numerical solution. *Mathematical Programming*, 199:1033–1106, 2023.
- [63] E. Tanguy, J. Delon, and N. Gozlan. Computing barycentres of measures for generic transport costs. *Preprint arXiv:2501.04016*, 2024.

- [64] B. Taşkesen, M.-C. Yue, J. Blanchet, D. Kuhn, and V. A. Nguyen. Sequential domain adaptation by synthesizing distributionally robust experts. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pages 10162–10172. PMLR, 2021.
- [65] B. Taşkesen, S. Shafieezadeh-Abadeh, D. Kuhn, and K. Natarajan. Discrete optimal transport with independent marginals is #P-hard. *SIAM Journal on Optimization*, 33(2):589–614, 2023.
- [66] A. B. Taylor. *Convex interpolation and performance estimation of first-order methods for convex optimization.* PhD thesis, Catholic University of Louvain, Louvain-la-Neuve, Belgium, 2017.
- [67] A. Vacher, B. Muzellec, A. Rudi, F. Bach, and F.-X. Vialard. A dimension-free computational upperbound for smooth optimal transport estimation. In *Proceedings of 34th Conference on Learning Theory*, volume 134, pages 4143–4173. PMLR, 2021.
- [68] C. Villani. *Topics in optimal transportation*, volume 58 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, 2003.
- [69] C. Villani. Optimal transport: Old and new, volume 338 of Grundlehren der mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]. Springer-Verlag, Berlin, 2009.
- [70] J. von Lindheim. Approximative algorithms for multi-marginal optimal transport and free-support Wasserstein barycenters. *Preprint arXiv:2202.00954*, 2022.

AREA OF DECISION SCIENCES, INSEAD, 1 AYER RAJAH AVE, 138676 SINGAPORE *Email address*: zeyi.chen@insead.edu

DIVISION OF MATHEMATICAL SCIENCES, NANYANG TECHNOLOGICAL UNIVERSITY, 21 NANYANG LINK, 637371 SINGA-PORE

Email address: ariel.neufeld@ntu.edu.sg

DIVISION OF MATHEMATICAL SCIENCES, NANYANG TECHNOLOGICAL UNIVERSITY, 21 NANYANG LINK, 637371 SINGA-PORE

Email address: qikun.xiang@ntu.edu.sg