
Sliced Multi-Marginal Optimal Transport

Samuel Cohen¹
Brandon Amos⁴

Alexander Terenin^{2, 5}
Marc Peter Deisenroth¹

Yannik Pitcan³
K S Sesh Kumar⁵

¹Centre for Artificial Intelligence, University College London
³University of California, Berkeley

⁴Facebook AI Research

²University of Cambridge
⁵Imperial College London

Abstract

Multi-marginal optimal transport enables one to compare multiple probability measures, which increasingly finds application in multi-task learning problems. One practical limitation of multi-marginal transport is computational scalability in the number of measures, samples and dimensionality. In this work, we propose a multi-marginal optimal transport paradigm based on random one-dimensional projections, whose (generalized) distance we term the *sliced multi-marginal Wasserstein distance*. To construct this distance, we introduce a characterization of the one-dimensional multi-marginal Kantorovich problem and use it to highlight a number of properties of the sliced multi-marginal Wasserstein distance. In particular, we show that (i) the sliced multi-marginal Wasserstein distance is a (generalized) metric that induces the same topology as the standard Wasserstein distance, (ii) it admits a dimension-free sample complexity, (iii) it is tightly connected with the problem of barycentric averaging under the sliced-Wasserstein metric. We conclude by illustrating the sliced multi-marginal Wasserstein on multi-task density estimation and multi-dynamics reinforcement learning problems.

1 Introduction

Optimal transport is a framework for defining meaningful metrics between probability measures [25, 31]. These metrics find a wide range of applications, such as generative modeling [11, 18], Bayesian inference [28], imitation learning [15], graph matching and averaging [32, 33]. Multi-marginal optimal transport [17] studies ways of comparing more than two probability measures in a geometrically meaningful way. Multi-marginal distances defined using this paradigm are often useful in settings where sharing geometric structure is useful, such as multi-task learning. In particular, they have been applied for training multi-modal generative adversarial networks [12], clustering [7], and computing barycenters of measures [4].

Following the establishment of key theoretical results, including by Agueh and Carlier [1], Gangbo and Świąch [17], and Pass [24], research is shifting toward applications. This motivates a need for practical algorithms for the multi-marginal setting [20]. Standard approaches based on linear programming and entropic regularization scale exponentially with the number of measures, and/or the dimension of the space [6, 29]. A number of recent works have therefore studied settings, where multi-marginal transport problems can be efficiently solved via low-rank structures on the underlying cost function [4], but exponential cost in the dimension remains [2, 3].

In parallel, a number of works on *sliced transport* [9] developed techniques for scalable transport, which (i) derive a closed form for a problem in a single dimension, and (ii) extend it into higher dimensions via random linear projections (slicing) and thereby inherit the complexity of the one-dimensional problem. This strategy has been shown effective in the classical Wasserstein [8, 9, 16,

19, 23, 27] and Gromov–Wasserstein [30] settings between pairs of measures, but has not yet been applied to settings with more than two measures.

In this paper, we address this gap and propose *sliced multi-marginal transport*, providing a scalable analog of the multi-marginal Wasserstein distance. To do so, we derive a closed-form expression for multi-marginal Wasserstein transport in one dimension, which lifts to a higher-dimensional analog via slicing. This one-dimensional closed-form expression can be computed with a complexity of $\mathcal{O}(PN \log N)$, where P is the number of measures and N is the number of samples per measure. Sliced multi-marginal Wasserstein (\mathcal{SMW}) can be estimated by Monte Carlo in $\mathcal{O}(KPN \log N)$, where K is the number of Monte Carlo samples.

Furthermore, we study \mathcal{SMW} 's theoretical properties. We prove that (i) it is a generalized metric, whose associated topology is the topology of weak convergence, (ii) its sample complexity is dimension free, just like the sliced Wasserstein case involving two measures, and (iii) sliced multi-marginal transport is closely connected with the problem of barycentric averaging under the sliced Wasserstein metric. We also showcase applications, where we focus on multi-task learning on probability spaces, where sharing knowledge across tasks can be beneficial and sliced multi-marginal Wasserstein can be used as a regularizer between task-specific models.

2 Background

Multi-marginal optimal transport [17] is a class of optimization problems for comparing multiple measures $\mu_1, \dots, \mu_P \in \mathcal{M}(\mathbb{R}^d)$, all supported on the metric space $(\mathbb{R}^d, \|\cdot\|_2)$. The most common such problem is computing the multi-marginal Wasserstein distance, defined as

$$\mathcal{MW}^2(\mu_1, \dots, \mu_P) = \min_{\pi \in \Pi(\mu_1, \dots, \mu_P)} \int_{(\mathbb{R}^d)^P} c(\mathbf{x}_1, \dots, \mathbf{x}_P) d\pi(\mathbf{x}_1, \dots, \mathbf{x}_P), \quad (1)$$

where $c : \mathbb{R}^d \times \dots \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a cost function and $\Pi(\mu_1, \dots, \mu_P)$ is the set of probability measures in $\mathcal{M}((\mathbb{R}^d)^P)$ with marginals μ_1, \dots, μ_P . We focus on the barycentric cost of Agueh and Carlier [1] and Gangbo and Świąch [17], given by

$$c(\mathbf{x}_1, \dots, \mathbf{x}_P) = \sum_{p=1}^P \beta_p \left\| \mathbf{x}_p - \sum_{j=1}^P \beta_j \mathbf{x}_j \right\|^2, \quad \beta_1, \dots, \beta_P \geq 0, \quad \sum_{p=1}^P \beta_p = 1. \quad (2)$$

This cost was originally motivated from an economics-inspired perspective, but is also often preferable because it leads to connections with barycentric averaging [1], giving it a simple interpretation. It also recovers the Wasserstein distance with squared 2-Euclidean cost in the case $P = 2$ (up to constants), referred to as \mathcal{W} . Algorithms for estimating (1) from a set of samples scale exponentially with the number of measures P and/or the dimension d of the ground space [2, 4, 6].

\mathcal{MW} is useful in multi-task settings for regularizing measures μ_1, \dots, μ_P by adding $\mathcal{MW}(\mu_1, \dots, \mu_P)$ to a multi-task loss. It can also be used in a setting, where we aim for a model output μ to be close to a given set of measures ν_1, \dots, ν_P , which can be done by introducing a loss of the form $\mathcal{MW}(\mu, \nu_1, \dots, \nu_P)$ and minimizing it with respect to μ .

Sliced transport. With the usual Euclidean-type cost structures, the Wasserstein distance between pairs of one-dimensional discrete measures can be computed efficiently using *sorting* with $\mathcal{O}(N \log N)$ complexity. More generally, we can consider the average distance between measures projected onto \mathbb{R} along random axis, which gives [8, 9]

$$\mathcal{SW}^2(\mu, \nu) = \int_{S_{d-1}} \mathcal{W}^2(M_{\#}^{\theta}(\mu), M_{\#}^{\theta}(\nu)) d\Theta(\theta), \quad (3)$$

where $M^{\theta}(\mathbf{x}) = \mathbf{x}^T \theta$, $(\cdot)_{\#}$ denotes the push-forward of measures, and Θ is the uniform distribution on the unit sphere S_{d-1} . We sample from $M_{\#}^{\theta}(\mu)$ by sampling from μ and projecting onto θ .

A fundamental result by Bonnotte [9] is that \mathcal{SW} is a metric that metrizes the topology of weak convergence—the *exact same* topology as \mathcal{W} . \mathcal{SW} can be estimated via Monte Carlo and preserves the computational complexity of estimating \mathcal{W} on \mathbb{R} , which is $\mathcal{O}(N \log N)$. Owing to the Monte Carlo nature, the sample complexity of \mathcal{SW} is dimension free [9, 22], in contrast with the exponential

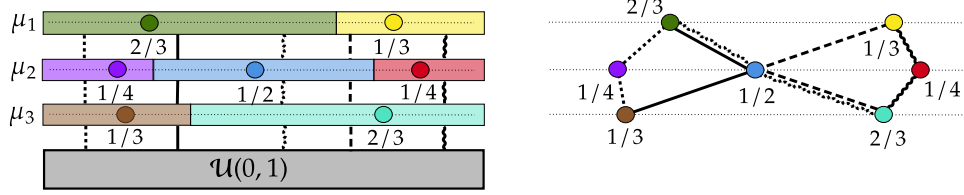


Figure 1: Illustration of the optimal coupling's structure on \mathbb{R} between discrete measures μ_1, μ_2 and μ_3 . Points are samples of each measures, with weights next to them. Left: histogram of measures (horizontal); joint samples are obtained by sampling a (black) line uniformly (drawn vertically), and picking points that are associated with the bin intersected by that line. Right: Corresponding triples of points that are aligned according to the coupling are linked by a pair of lines.

dependency of the Wasserstein distance on dimension. The combination of good computational and statistical properties makes \mathcal{SW} an attractive choice for minimization problems on measure spaces, including generative modeling and imitation learning [15, 16]. This immediately raises the question whether \mathcal{SW} extends to the multi-marginal case so that it preserves its key appealing properties.

3 Sliced Multi-Marginal Transport

To proceed toward a suitable notion of sliced multi-marginal optimal transport, we begin by developing a probabilistic analogy to understand the coupling structure that arises in one-dimensional transport when considering multiple measures. This enables us to derive suitably-closed-form expressions from which sliced multi-marginal Wasserstein distances can be built.

3.1 One-dimensional Multi-Marginal Transport

In optimal transport, couplings between probability measures form one of the standard objects of study. One way to understand the structure of a coupling is to introduce a set of random variables $y_i : \Omega \rightarrow \mathbb{R}$ on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ whose joint distribution is the coupling of interest. Consider the one-dimensional Wasserstein formula

$$\mathcal{W}^2(\mu_1, \mu_2) = \int_0^1 |C_{\mu_1}^{-1}(x) - C_{\mu_2}^{-1}(x)|^2 dx, \quad (4)$$

where $C_{\mu_1}^{-1}, C_{\mu_2}^{-1}$ are the generalized quantile functions of μ_1, μ_2 . If we define $y_1 = C_{\mu_1}^{-1}(x)$ and $y_2 = C_{\mu_2}^{-1}(x)$, taking $([0, 1], \mathcal{B}(0, 1), \mathcal{U}(0, 1))$ as our probability space, we can write (4) as

$$\mathcal{W}^2(\mu_1, \mu_2) = \mathbb{E}_{y_1, y_2 \sim \Pi} [|y_1 - y_2|^2] \quad \Pi = (C_{\mu_1}^{-1}, C_{\mu_2}^{-1})_{\#} \mathcal{U}(0, 1). \quad (5)$$

This reveals that the optimal coupling admits a very specific structure: it is the pushforward measure induced by an underlying uniform random variable. The one-dimensional Wasserstein distance is thus simply the average squared distance of a pair of random variables y_1 and y_2 , where (a) we sample both y_1 and y_2 by the generalized quantile method, and (b) we *share the underlying uniform random numbers* used in the sampling. We prove that this view is general and extends to the multi-marginal case, even in the case of the more elaborate cost structure introduced in Section 2.

Proposition 1. *If $\mu_1, \dots, \mu_P \in \mathcal{M}(\mathbb{R})$ and $\mathcal{U}(0, 1)$ is the uniform measure, then*

$$\mathcal{MW}^2(\mu_1, \dots, \mu_P) = \int_0^1 \sum_{p=1}^P \beta_p \left| C_{\mu_p}^{-1}(x) - \sum_{j=1}^P \beta_j C_{\mu_j}^{-1}(x) \right|^2 dx, \quad (6)$$

and the optimal coupling solving (1) is of the form

$$\pi^* = (C_{\mu_1}^{-1}, \dots, C_{\mu_P}^{-1})_{\#} \mathcal{U}(0, 1). \quad (7)$$

Proposition 1 shows the optimal coupling is the push-forward of a uniform distribution through the generalized quantiles of each measure. Obtaining joint samples from the coupling can hence be

done by sampling from the uniform distribution and mapping through each quantile function. This extends the result by Carlier et al. [13] to the setting where absolute continuity is not assumed. In the discrete case, we can simplify this further by introducing the sorting idea used in the one-dimensional Wasserstein case, to deduce the following.

Corollary 2. *If measures $\mu_1, \dots, \mu_P \in \mathcal{M}(\mathbb{R})$ are discrete and uniform with N atoms, i.e., $\mu_p = \frac{1}{N} \sum_{i=1}^N \delta_{\tilde{x}_i^{(p)}}$, with $\tilde{x}_1^{(p)} \leq \dots \leq \tilde{x}_N^{(p)}$, for $p = 1, \dots, P$, then*

$$\mathcal{MW}^2(\mu_1, \dots, \mu_P) = \frac{1}{N} \sum_{i,p=1}^{N,P} \beta_p \left| \tilde{x}_i^{(p)} - \sum_{j=1}^P \beta_j \tilde{x}_i^{(j)} \right|^2. \quad (8)$$

In particular, this means that the complexity of computing the multi-marginal Wasserstein in one dimension in the discrete uniform case is $\mathcal{O}(PN \log N)$ —the cost of sorting. This establishes the necessary results in one dimension, and we generalize them to the higher-dimensional case via slicing.

3.2 Sliced Multi-Marginal Wasserstein Distance

To define the sliced multi-marginal Wasserstein distance, we average the expressions given in (6) along one-dimensional random projections, which gives

$$\mathcal{SMW}^2(\mu_1, \dots, \mu_P) = \int_{S_{d-1}} \int_0^1 \sum_{p=1}^P \beta_p \left| C_{\mu_p^\theta}^{-1}(x) - \sum_{j=1}^P \beta_j C_{\mu_j^\theta}^{-1}(x) \right|^2 dx d\Theta(\theta), \quad (9)$$

where $\mu_j^\theta = M_{\theta\#}(\mu_j)$ for $j = 1, \dots, P$. \mathcal{SMW} in (9) can be estimated via Monte Carlo in $\mathcal{O}(KPN \log N)$, where K is the number of Monte Carlo samples (projections).

Topological properties We now study \mathcal{SMW} 's topological properties. We first show that \mathcal{SMW} is the weighted mean of sliced Wasserstein distances between pairs of measures.

Proposition 3. *Let $\mu_1, \dots, \mu_P \in \mathcal{M}(\mathbb{R}^d)$. We have that*

$$\mathcal{SMW}^2(\mu_1, \dots, \mu_P) = \frac{1}{2} \sum_{i,j=1}^P \beta_i \beta_j \mathcal{SW}^2(\mu_i, \mu_j). \quad (10)$$

Proposition 3 is useful in deriving statistical and topological properties of \mathcal{SMW} . It is however more efficient to estimate it via our closed-form formula for multi-marginal transport – see (9). This leads to a computational complexity of $\mathcal{O}(KPN \log N)$, whereas naively implementing (10) scales in $\mathcal{O}(KP^2N \log N)$. Furthermore, as the sliced-Wasserstein metric is upper-bounded by the Wasserstein [9], an immediate consequence of Proposition 3 is that

$$\mathcal{SMW}^2(\mu_1, \dots, \mu_P) \stackrel{(10)}{=} \frac{1}{2} \sum_{i,j=1}^P \beta_i \beta_j \mathcal{SW}^2(\mu_i, \mu_j) \leq \frac{1}{2} \sum_{i,j=1}^P \beta_i \beta_j \mathcal{W}^2(\mu_i, \mu_j). \quad (11)$$

A reverse inequality also follows directly (see corollary 12), which shows that \mathcal{SMW} gives rise to the topology of weak convergence—one of the key properties that made \mathcal{SW} an attractive choice in the first place. We now study metric properties of \mathcal{SMW} .

Proposition 4. *\mathcal{SMW} is a generalized metric.*

In particular, this means that \mathcal{SMW} is (i) non-negative, (ii) zero if and only if all measures are identical, (iii) permutation-equivariant, and (iv) satisfies a generalized triangle inequality involving multiple measures. Hence, \mathcal{SMW} is well-behaved topologically-wise as it is a generalized metric inducing weak convergence. We continue by studying \mathcal{SMW} 's statistical properties.

Statistical Properties In the following proposition, we assess the impact of the number of samples and random projections used to estimate \mathcal{SMW} .

Proposition 5. If $\mu_1, \dots, \mu_P \in \mathcal{M}(\mathbb{R}^d)$, and assuming \mathcal{W}^2 has sample complexity $\rho(N)$ on \mathbb{R} , then,

$$\mathbb{E}[\mathcal{SMW}^2(\mu_1, \dots, \mu_P) - \mathcal{SMW}^2(\hat{\mu}_1, \dots, \hat{\mu}_P)]^2 \leq \frac{1}{2}\rho(N), \quad (12)$$

where $\hat{\mu}_p$ refers to empirical measures with N samples.

Proposition 5 shows that the sample complexity of \mathcal{SMW} is dimension-free—this stands in contrast to the sample complexity of the multi-marginal Wasserstein, which is exponential in the dimension. In practice, we use Monte Carlo sampling to compute \mathcal{SMW} , which introduces additional error. To understand this error, we examine \mathcal{SMW} 's projection complexity.

Proposition 6. Let $\mu_1, \dots, \mu_P \in \mathcal{M}(\mathbb{R}^d)$, and define $\overline{\mathcal{SMW}}$ the approximation obtained by uniformly picking L projections on S_{d-1} , then

$$\mathbb{E} \left[\overline{\mathcal{SMW}}^2(\mu_1, \dots, \mu_P) - \mathcal{SMW}^2(\mu_1, \dots, \mu_P) \right]^2 \leq L^{-1/2} \text{Var}_{\theta} \left[\mathcal{MW}^2(\mu_1^{\theta}, \dots, \mu_P^{\theta}) \right], \quad (13)$$

where θ follows the uniform distribution on S_{d-1} and $\mu_p^{\theta} = M_{\#}^{\theta}(\mu_p)$.

This shows that the quality of Monte Carlo estimates of \mathcal{SMW} is controlled by number of projections and the variance of evaluations of the base multi-marginal Wasserstein in 1D.

Connection to Barycenters We now study connections of \mathcal{SMW} to the problem of barycentric averaging, which extends the notion of a *mean* to more general settings. Let $\mathcal{D} : \mathcal{M}(\mathbb{R}^d) \times \mathcal{M}(\mathbb{R}^d) \rightarrow \mathbb{R}$ be a discrepancy on the space of probability measures. Recall that the *barycenter* of P measures μ_1, \dots, μ_P is defined as

$$\mu^* = \arg \min_{\mu \in \mathcal{M}(\mathbb{R}^d)} \mathcal{F}(\mu), \quad \mathcal{F}(\mu) = \sum_{p=1}^P \mathcal{D}(\mu_p, \mu). \quad (14)$$

Barycentric averaging is well-studied from theoretical and computational view-points, notably under the squared Wasserstein [14], sliced Wasserstein [8] and Gromov–Wasserstein [26] metrics.

Proposition 7. Let $\mu_1, \dots, \mu_P \in \mathcal{M}(\mathbb{R}^d)$, $\sum_{p=1}^P \beta_p = 1$. Furthermore, let $\hat{\beta}_p$ be augmented multi-marginal weights, so that for $m \in [0, 1]$ it holds that $\hat{\beta}_p = m\beta_p$ for $p = 1, \dots, P$, $\sum_{p=1}^{P+1} \hat{\beta}_p = 1$, and $\mathcal{D} = \mathcal{SW}^2$. Then

$$\arg \min_{\mu \in \mathcal{M}(\mathbb{R}^d)} \mathcal{SMW}^2(\mu_1, \dots, \mu_P, \mu) = \arg \min_{\mu \in \mathcal{M}(\mathbb{R}^d)} \mathcal{F}(\mu), \quad (15)$$

where β is the weight vector of \mathcal{F} and $\hat{\beta}$ is the weight vector of \mathcal{SMW} .

Proposition 7 reveals a connection between sliced multi-marginal transport and barycenters under the sliced-Wasserstein: the measure that is closest to μ_1, \dots, μ_P in \mathcal{SMW} is actually the barycenter of such measures under \mathcal{SW} . We continue by studying smoothness of \mathcal{SMW} as a loss function.

Differentiability Sliced Wasserstein variants are desirable candidate losses for learning on probability spaces thanks to their smoothness properties. We show \mathcal{SMW} inherits these properties.

Proposition 8. Let $\mu_1, \dots, \mu_P \in \mathcal{M}(\mathbb{R}^d)$ be discrete measures with N atoms, which we gather into matrices $\{\mathbf{X}^{(p)}\}_{p=1}^P$, and similarly define $\mu_{\mathbf{X}}$ with atoms \mathbf{X} . Assume \mathbf{X} has distinct points. Then \mathcal{SMW}^2 is smooth with gradient

$$\nabla_{\mathbf{X}} \mathcal{SMW}^2(\mu_1, \dots, \mu_P, \mu_{\mathbf{X}}) = \beta_{P+1} \sum_{p=1}^P \beta_p \int_{S_{d-1}} \mathbf{X}_{\theta} - (\mathbf{X}_{\theta}^{(p)} \circ \sigma_{\mathbf{X}_{\theta}} \circ \sigma_{\mathbf{X}_{\theta}^{(p)}}^{-1}) \, d\Theta(\theta), \quad (16)$$

where $\sigma_{\mathbf{X}}$ is the permutation that sorts atoms of \mathbf{X} .

Proposition 8 shows that \mathcal{SMW}^2 is smooth almost everywhere, and is hence well-suited for multi-task learning, as it allows to compare multiple task-representative probability measures. We illustrate this in Figure 2. Here, we consider the problem $\min_{\mu} \mathcal{SMW}^2(\mu, \nu_1, \dots, \nu_4)$, amounting to estimating the sliced barycenter of μ_1, \dots, μ_4 (see Proposition 7), and solve it iteratively via the gradient flow $\partial \mu_t = -\nabla \mathcal{SMW}^2(\mu_t, \nu_1, \dots, \nu_P)$, following Bonneel et al. [8] in the pairwise case.

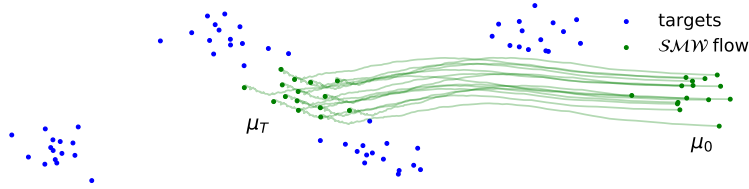


Figure 2: Gradient flow $\partial\mu_t = -\nabla\mathcal{SMW}^2(\mu_t, \nu_1, \dots, \nu_P)$ starting from a randomly initialized Gaussian μ_0 . It is solved iteratively following Bonneel et al. [8].

4 Multi-Task Learning with Sliced Multi-marginal Optimal Transport

In the previous section, we proposed a multi-marginal metric between probability measures, which avoids exponential computational and statistical complexities and is thus practical for applications where a large number of samples N , number of measures P , or dimension d is of interest. \mathcal{SMW} allows us to evaluate the closeness of probability measures μ_1, \dots, μ_P , which makes it a good candidate regularizer in multi-task learning settings over probability spaces, by encouraging shared global structure across tasks through closeness in sliced multi-marginal geometry. We now sketch potential areas of applications of \mathcal{SMW} in the context of multi-task learning on spaces of probability measures, and illustrate examples in density estimation and multi-dynamics reinforcement learning.

4.1 Density Estimation with Shared Structure

Consider P target measures μ_1, \dots, μ_P , which we aim to approximate by parametric models ν_1, \dots, ν_P , such as for instance generative adversarial networks. In applications, it is often the case that these measures are affected by issues related to *distributional shift* [5], which prevents us from obtaining accurate empirical samples of μ_1, \dots, μ_P . One way to counteract these issues is to introduce a shared structure between the measures, which can be enforced through \mathcal{SMW} regularization.

For example, consider empirical estimates $\hat{\mu}_1, \dots, \hat{\mu}_P$ of μ_1, \dots, μ_P , which are corrupted because no data is available in certain regions of each measure’s support. Here, reconstruction of μ_1, \dots, μ_P is only possible through the use of shared structure on the generative models ν_1, \dots, ν_P , which we can enforce by using $\mathcal{SMW}(\nu_1, \dots, \nu_P)$ as a regularizer. This results in the optimization problem

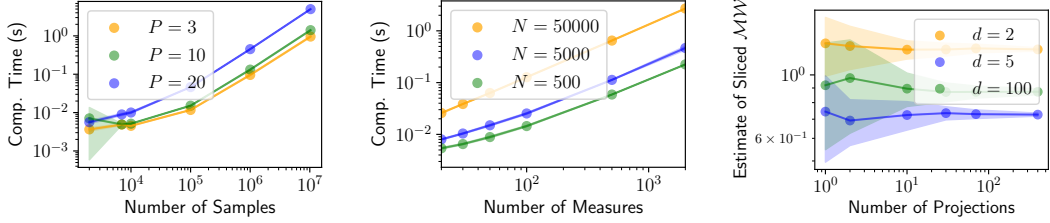
$$\arg \min_{\nu_1, \dots, \nu_P} \sum_{p=1}^P \underbrace{\mathcal{SW}^2(\mu_p, \nu_p)}_{\text{local loss}} + \gamma \underbrace{\mathcal{SMW}^2(\nu_1, \dots, \nu_P)}_{\text{global loss (shared)}}, \quad (17)$$

where $\mathcal{SW}^2(\mu_p, \nu_p)$ ensures that the respective generative models $(\nu_p)_{p=1}^P$ approximates targets $(\mu_p)_{p=1}^P$, and $\mathcal{SMW}^2(\nu_1, \dots, \nu_P)$ ensures shared structure is present in the loss.

4.2 Multi-Dynamics Reinforcement Learning with Shared Structure

We now consider the problem of reinforcement learning in settings where the dynamics change. In order to speed up learning, we use \mathcal{SMW} to share structure across different environments in this multi-dynamics reinforcement learning problem. Sharing knowledge is not only useful to bias (and thereby speed up) learning, but it is also useful in settings, where agents are ill informed, e.g., due to sparse reward signals. With a shared structure, these agents can learn from other agents. Here, the challenge is in effectively utilizing information from other agents in spite of differences in their respective environments. In the following, we focus on this setting.

Consider P identical-task agents in finite-horizon Markov decision processes $(\mathcal{S}, \mathcal{A}, \mathcal{T}_p, r_p^{\text{env}})$, where \mathcal{S} is the state space and \mathcal{A} is the action space, both shared by all agents, $T_p(\mathbf{x}_t^{(p)}, \mathbf{a}_t^{(p)}) = \mathbf{x}_{t+1}^{(p)}$ is the transition model of agent p , which varies across agents, and r_p^{env} is the environment’s reward function. Since different agents’ tasks are identical, sharing structure can be beneficial. We consider the case, where some agents receive rewards $r_p^{\text{env}} = 0$. These agents are *uninformed* and can only learn via a shared structure that allows to transfer knowledge from other agents. Structure sharing is done by



(a) Computational time (log-log scale, mean \pm standard deviation over 5 runs) for computing the sliced multi-marginal distance in seconds against the number of samples for various P . (b) Computational time (log-log scale, mean \pm standard deviation over 5 runs) for computing the sliced multi-marginal distance in seconds against the number of measures, $d = 10$ for various N . (c) Mean SMW (\pm standard deviations) sliced multi-marginal distance against the number of projections for $P = 5$ measures with $N = 250$ samples.

Figure 3: Properties of the sliced multi-marginal distance. (a) computational time as a function of the number of samples; (b) computational time as a function of the number of measures; (c) accuracy as a function of the number of projections

augmenting the agent-specific reward function with a global multi-task reward term. In particular, define the augmented reward R_p as

$$R_p(\mathbf{x}_t^{(p)}) = \underbrace{r_p^{\text{env}}(\mathbf{x}_t^{(p)})}_{\text{agent specific (local)}} + \gamma \underbrace{r^{\text{mul}}(\mathbf{x}_t^{(p)}, \mathbf{X})}_{\text{multi-task reward (shared/global)}}, \quad (18)$$

where $\mathbf{X} = \{\mathbf{x}_t^{(p)}\}_{p,t=1}^{P,T}$ is the collection of all states of every agent at all time steps, $r_p^{\text{env}}(\mathbf{x}_t^{(p)})$ is the single-task reward of the p^{th} environment and $r^{\text{mul}}(\mathbf{x}_t^{(p)}, \mathbf{X})$ is a (multi-task) reward signal. The former provides task-specific information about the task to be solved by agent p , while the latter allows for agents to share structure through the history of their state trajectories. If $r_p^{\text{env}} = 0$ for a given agent, then this agent can only learn through the shared structure arising from the shared reward r^{mul} . Finally, γ is a regularizer that controls the influence of shared structure on the overall learning.

We now describe the shared reward r^{mul} . Denote $\mu_p = \frac{1}{T} \sum_{t=1}^T \delta_{\mathbf{x}_t^{(p)}}$, which allows us to interpret the rollout of agent p as a discrete probability measure supported on the state space. Then,

$$r^{\text{mul}}(\mathbf{x}_t^{(p)}, \mathbf{X}) = -\frac{\beta_p}{K} \sum_{k=1}^K \left| \langle \mathbf{x}_t^{(p)} - \sum_{j=1}^P \beta_j \mathbf{x}_{\eta_{p,j,k}(t)}^{(j)}, \boldsymbol{\theta}_k \rangle \right|^2, \quad (19)$$

where $\eta_{p,j,k}$ returns the index of the atom in μ_j that is aligned with state $\mathbf{x}_t^{(p)}$ after projecting on (Monte Carlo-sampled) $(\boldsymbol{\theta}_k)_{k=1}^K$ and sorting all projected states. Intuitively, the reward signal attributed to the state $\mathbf{x}_t^{(p)}$ of agent p at time t is computed by projecting all measures onto K vectors, gathering all states that are aligned with $\mathbf{x}_t^{(p)}$ for each projection $\boldsymbol{\theta}_k$, and summing squared distances between them.

Remark. *The barycentric cost structure with non-uniform weights β is particularly attractive in this setting, as it allows to give more weight to the communication arising from agents that perform well in their own environment. For instance, we can use Boltzmann weights*

$$\beta_p \propto \exp \left(\alpha \sum_{t=1}^T r_p^{\text{env}}(\mathbf{x}_t^{(p)}) \right), \quad (20)$$

where α is a temperature. It gives more weight in the reward to agents performing best.

We train all agents simultaneously by maximizing

$$\mathbb{E}_{\pi_1, \dots, \pi_P} \left[\sum_{p=1}^P \sum_{t=1}^T R_p(\mathbf{x}_t^{(p)}) \right] = \mathbb{E}_{\pi_1, \dots, \pi_P} \left[\sum_{p=1}^P \sum_{t=1}^T \underbrace{r_p^{\text{env}}(\mathbf{x}_t^{(p)}) - \gamma SMW^2(\mu_1, \dots, \mu_P)}_{=R_p(\mathbf{x}_t^{(p)})} \right] \quad (21)$$

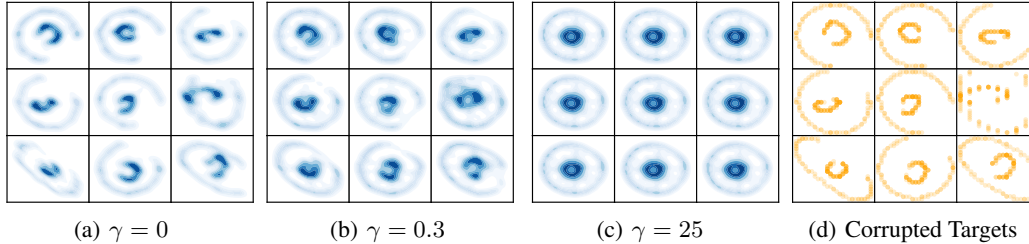


Figure 4: Multi-task density estimation experiment applied on corrupted nested ellipses (plotted in orange), using \mathcal{SW}^2 as pairwise loss and \mathcal{SMW}^2 as regularizer. Learned models are plotted in blue. We use regularization coefficients $\gamma = 0$ in (a), $\gamma = 0.3$ in (b), $\gamma = 25$ in (c).

with respect to the parameters of policies $\pi_p, p = 1, \dots, P$. Note that the extra term in the augmented reward regularizes the objective via the sliced multi-marginal Wasserstein distance. \mathcal{SMW} thus enforces closeness of agents’ trajectories which allows to share structure across agents.

5 Experiments

We now illustrate the behavior of sliced multi-marginal transport in simple multi-task learning setups.

5.1 Scalability

Number of Samples (N). We study the impact of the number of samples on the computational time to compute the sliced multi-marginal distance in (9). In particular, we compute \mathcal{SMW} between $P = 3, 10, 20$ measures in \mathbb{R}^{10} , $\mu_p \sim \mathcal{N}(\mathbf{m}_p, \eta^2 \mathbf{I})$, where $p = 1, \dots, P$ for a fixed number of projections $K = 10$. Figure 3(a) shows the $\mathcal{O}(N \log N)$ scaling of \mathcal{SMW} . This enables computation of multi-marginal distances with over 10^7 samples and a large number of measures.

Number of Measures (P). We now examine scaling with respect to the number of measures P . Figure 3(b) shows the time required to compute \mathcal{SMW} against $N = 500, 5000, 50000$ measures. We observe the expected linear scaling of \mathcal{SMW} .

Number of Projections (K). Finally, we consider the impact of the number of projections on the estimation of \mathcal{SMW} for dimensions $d = 2, 5, 20$. We set $N = 250$, and $P = 5$. Monte Carlo estimation is used to estimate \mathcal{SMW} . Figure 3(c) shows the expected variance shrinkage as the number of projection grows, while the estimated mean converges to \mathcal{SMW} with rate $\mathcal{O}(\frac{1}{\sqrt{K}})$ and constant factors depending on dimension.

5.2 Multi-Task Density Estimation

We consider the multi-task density estimation setting of Section 4.1. Each target measures consist of a nested ellipse with corrupted samples. In particular, parts of each individual ellipse have been removed from each measure’s support. Using the multi-task learning setup allows for sharing knowledge of the structure of the target tasks across problems—namely, that all target measures have the overall shape of nested ellipses. Figures 4(a)–4(c) show the models obtained by multi-task training with regularization coefficients $\gamma = 0, 0.3, 25$. When $\gamma = 0$, measures are learned individually without any structure sharing. ν_1, \dots, ν_P hence collapse to the corrupted measures μ_1, \dots, μ_P . When structure is introduced ($\gamma > 0$) knowledge of the inherent nested ellipse structure is shared across tasks, which leads to solutions that have such structure (holes are filled), but that still preserve the task-specific orientations and ellipse width/height as long as the structure coefficient η is not too large. The latter causes the learned measures to be too close to each other. These effects can be seen in Figure 4(c). When this happens, all learned measures collapse to the barycenter.

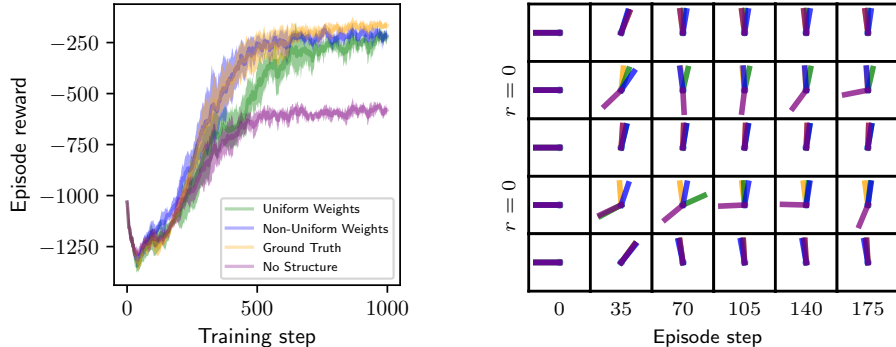


Figure 5: Multi-task ($P = 5$) RL experiment. Environments have different dynamics (different gravities), and $2/5$ agents have no environmental reward. Without shared structure, these agents do not solve their respective tasks (orange). By contrast, with shared structure, all agents learn accurate policies (green, blue), on par with agents trained without corrupted rewards (blue). Left: training curves (mean \pm standard deviation averaged over 5 runs), Right: states of agents for each task at the end of training (left to right refers to time t from 0 to 200).

5.3 Multi-Dynamics Reinforcement Learning

We consider a multi-task RL application in the setting of Section 4.2. In particular, we consider $P = 5$ pendulum swing-up tasks with different dynamics (gravities $g \in \{8, 9, 10, 11, 12\} \text{ m/s}^2$). States consist of angle and angular velocities, and actions of are torques. Environment rewards are dense as implemented in OpenAI Gym [10], and following Dadashi et al. [15], we transform the shared reward r^{mul} via $f(y) = e^{-5y}$. Two out of five agents do not receive any reward. All other agents share the same reward function. We consider agents trained with and without \mathcal{SMW} -based regularization, and consider the uniform and non-uniform barycentric weights β ; see Section 4.2 for more details. To facilitate learning, we lower-bound the weights of agents without reward. Policies are learned using Q -learning with function approximation on state observations.

Figure 5 shows the results. Training without regularization ($\gamma = 0$, blue curve) does not allow the two agents without environment rewards ($r_p^{\text{env}} = 0$) to solve their respective tasks. By contrast, with regularization, all agents (even those with no environment reward) solve their respective tasks (green, blue) as well as if all agents were receiving environmental rewards (orange). Agents with non-uniform regularization significantly outperform agents with uniform weights, showing that giving more weight in the regularizer to stronger agents is helpful. Overall, this demonstrates that knowledge transfer via the shared reward structure can be effective. In particular, the regularization-based rewards encourage the state trajectories of all agents to be close under the sliced multi-marginal geometry. Hence, agents without environment rewards learn to *follow* agents trained with environment rewards. This is possible because of similarity of environments and of agent goals, so that agent rollouts share geometric structure.

6 Conclusion

In this work, we proposed a scalable multi-marginal optimal transport distance. Our main idea is to derive a closed-form formula for multi-marginal optimal transport in 1D in the general case and to extend it into a higher-dimensional metric via slicing. We show it is well-behaved topologically, and in particular that it is a generalized metric. We also show it is well-behaved statistically with dimension-free sample complexity (modulo a caveat arising from projection complexity). We derive a range of other results illustrating the simple and intuitive geometric structure of sliced multi-marginal transport. Finally, we propose areas of applications of sliced multi-marginal transport in the context of multi-task learning on probability spaces, and concrete instantiations in density estimation, and reinforcement learning. We hope these contributions enable practitioners in reinforcement learning, generative modeling and other areas to share structure across tasks in a geometrically-motivated way. Our work relies on the assumption that tasks live on the same space, and share structure. Future work

extends our approach to allow for multi-task learning on incomparable spaces, enabling structure sharing in more general set-ups, for instance via Gromov–Wasserstein-like techniques.

Acknowledgments

SC was supported by the Engineering and Physical Sciences Research Council (grant number EP/S021566/1).

References

- [1] M. Agueh and G. Carlier. Barycenters in the Wasserstein Space. *SIAM Journal on Mathematical Analysis*, 43(2):904–924, 2011. Cited on pages 1, 2.
- [2] J. Altschuler and E. Boix-Adsera. Wasserstein Barycenters are NP-hard to Compute. *arXiv:2101.01100*, 2021. Cited on pages 1, 2.
- [3] J. M. Altschuler and E. Boix-Adserà. Hardness results for Multimarginal Optimal Transport problems. *arXiv:2012.05398*, 2020. Cited on page 1.
- [4] J. M. Altschuler and E. Boix-Adserà. Polynomial-time Algorithms for Multimarginal Optimal Transport Problems with Structure. *arXiv:2008.03006*, 2020. Cited on pages 1, 2.
- [5] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané. Concrete Problems in AI Safety. *arXiv:1606.06565*, 2016. Cited on page 6.
- [6] J.-D. Benamou, G. Carlier, M. Cuturi, L. Nenna, and G. Peyré. Iterative Bregman Projections for Regularized Transportation Problems. *SIAM Journal on Scientific Computing*, 37(2):A1111–A1138, 2015. Cited on pages 1, 2.
- [7] J. Bento and L. Mi. Multi-Marginal Optimal Transport Defines a Generalized Metric. *arXiv:2001.11114*, 2020. Cited on page 1.
- [8] N. Bonneel, J. Rabin, G. Peyré, and H. Pfister. Sliced and Radon Wasserstein Barycenters of Measures. *Journal of Mathematical Imaging and Vision*, 51(1):22–45, 2015. Cited on pages 1, 2, 5, 6, 20, 21.
- [9] N. Bonnotte. Unidimensional and Evolution Methods for Optimal Transportation, 2013. Cited on pages 1, 2, 4, 18.
- [10] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. OpenAI Gym. *arXiv:1606.01540*, 2016. MIT License. Cited on pages 9, 21.
- [11] C. Bunne, D. Alvarez-Melis, A. Krause, and S. Jegelka. Learning Generative Models across Incomparable Spaces. In *ICML*, 2019. Cited on page 1.
- [12] J. Cao, L. Mo, Y. Zhang, K. Jia, C. Shen, and M. Tan. Multi-marginal Wasserstein GAN. In *NeurIPS*, 2019. Cited on page 1.
- [13] G. Carlier, A. Oberman, and E. Oudet. Numerical Methods for Matching for Teams and Wasserstein Barycenters. *ESAIM*, 2015. Cited on page 4.
- [14] M. Cuturi and A. Doucet. Fast Computation of Wasserstein Barycenters. In *ICML*, 2014. Cited on page 5.
- [15] R. Dadashi, L. Hussenot, M. Geist, and O. Pietquin. Primal Wasserstein Imitation Learning. *arXiv:2006.04678*, 2020. Cited on pages 1, 3, 9, 21.
- [16] I. Deshpande, Y.-T. Hu, R. Sun, A. Pyrros, N. Siddiqui, S. Koyejo, Z. Zhao, D. A. Forsyth, and A. G. Schwing. Max-Sliced Wasserstein Distance and Its Use for GANs. In *CVPR*, 2019. Cited on pages 1, 3.
- [17] W. Gangbo and A. Świąch. Optimal maps for the multidimensional Monge-Kantorovich problem. *Communications on Pure and Applied Mathematics*, 51(1):23–45, 1998. Cited on pages 1, 2.
- [18] A. Genevay, G. Peyre, and M. Cuturi. Learning Generative Models with Sinkhorn Divergences. In *AISTATS*, 2018. Cited on page 1.
- [19] S. Kolouri, K. Nadjahi, U. Simsekli, R. Badeau, and G. Rohde. Generalized Sliced Wasserstein Distances. In *NeurIPS*. 2019. Cited on page 1.
- [20] T. Lin, N. Ho, M. Cuturi, and M. I. Jordan. On the Complexity of Approximating Multimarginal Optimal Transport. *arXiv:1910.00152*, 2019. Cited on page 1.

- [21] K. Nadjahi, A. Durmus, L. Chizat, S. Kolouri, S. Shahrampour, and U. Simsekli. Statistical And Topological Properties of Sliced Probability Divergences. *arXiv:2003.05783*, 2020. Cited on page 16.
- [22] K. Nadjahi, A. Durmus, L. Chizat, S. Kolouri, S. Shahrampour, and U. Simsekli. Statistical and Topological Properties of Sliced Probability Divergences. In *NeurIPS*, 2020. Cited on pages 2, 18, 19.
- [23] K. Nguyen, N. Ho, T. Pham, and H. Bui. Distributional Sliced-Wasserstein and Applications to Generative Modeling. In *ICLR*, 2021. Cited on page 2.
- [24] B. Pass. Multi-Marginal Optimal Transport: Theory and Applications. *arXiv:1406.0026*, 2014. Cited on page 1.
- [25] G. Peyré and M. Cuturi. Computational Optimal Transport. *Foundations and Trends in Machine Learning*, 2019. Cited on pages 1, 12, 16.
- [26] G. Peyré, M. Cuturi, and J. Solomon. Gromov-Wasserstein Averaging of Kernel and Distance Matrices. In *ICML*, 2016. Cited on page 5.
- [27] M. Rowland, J. Hron, Y. Tang, K. Choromanski, T. Sarlos, and A. Weller. Orthogonal Estimation of Wasserstein Distances. In *AISTATS*, 2019. Cited on page 2.
- [28] S. Srivastava, C. Li, and D. B. Dunson. Scalable Bayes via Barycenter in Wasserstein Space. *Journal of Machine Learning Research*, 19(1):312–346, Jan. 2018. Cited on page 1.
- [29] N. Tupitsa, P. Dvurechensky, A. Gasnikov, and C. A. Uribe. Multimarginal Optimal Transport by Accelerated Alternating Minimization. *CDC:6132–6137*, 2020. Cited on page 1.
- [30] T. Vayer, R. Flamary, N. Courty, R. Tavenard, and L. Chapel. Sliced Gromov-Wasserstein. In *NeurIPS*. 2019. Cited on page 2.
- [31] C. Villani. *Optimal Transport: Old and New*. Springer Science & Business Media, 2008. Cited on page 1.
- [32] H. Xu, D. Luo, and L. Carin. Scalable Gromov-Wasserstein Learning for Graph Partitioning and Matching. In *NeurIPS*. 2019. Cited on page 1.
- [33] H. Xu, D. Luo, H. Zha, and L. C. Duke. Gromov-Wasserstein Learning for Graph Matching and Node Embedding. In *ICML*, 2019. Cited on page 1.

A Proofs

A.1 Closed-form Formulas for Multimarginal Optimal Transport

For a measure $\mu \in \mathcal{M}(\mathbb{R})$, define its CDF $C_\mu : \mathbb{R} \rightarrow [0, 1]$ as

$$C_\mu(x) = \int_{-\infty}^x d\mu(y) \quad \forall x. \quad (22)$$

Also, define its pseudo-inverse $C_\mu^{-1} : [0, 1] \rightarrow \mathbb{R} \cup \{-\infty\}$ as

$$C_\mu^{-1}(r) = \min_x \{x \in \mathbb{R} \cup \{-\infty\} : C_\mu(x) \geq r\}. \quad (23)$$

This function is a generalization of the quantile function.

1D Multi-Marginal

Proposition 1. *If $\mu_1, \dots, \mu_P \in \mathcal{M}(\mathbb{R})$ and $\mathcal{U}(0, 1)$ is the uniform measure, then*

$$\mathcal{MW}^2(\mu_1, \dots, \mu_P) = \int_0^1 \sum_{p=1}^P \beta_p \left| C_{\mu_p}^{-1}(x) - \sum_{j=1}^P \beta_j C_{\mu_j}^{-1}(x) \right|^2 dx, \quad (24)$$

and the optimal coupling solving (1) is of the form

$$\pi^* = (C_{\mu_1}^{-1}, \dots, C_{\mu_P}^{-1})_{\#} \mathcal{U}(0, 1). \quad (25)$$

Proof. Our aim is to provide a closed form formula for

$$\mathcal{MW}^2(\mu_1, \dots, \mu_P) = \min_{\pi \in \Pi(\mu_1, \dots, \mu_P)} \int_{(\mathbb{R}^d)^P} \sum_{p=1}^P \beta_p |x_p - \sum_j \beta_j x_j|^2 d\pi(x_1, \dots, x_P), \quad (26)$$

where $\Pi(\mu_1, \dots, \mu_P)$ is the set of probability measures in $\mathcal{M}((\mathbb{R}^d)^P)$ with marginals μ_1, \dots, μ_P .

First, notice

$$\int_{(\mathbb{R}^d)^P} \sum_{p=1}^P \beta_p \|x_p - \sum_j \beta_j x_j\|^2 d\pi(x_1, \dots, x_P) \quad (27)$$

$$= \sum_{p=1}^P \beta_p \int_{\mathbb{R}^d} |x_p|^2 d\mu_p - 2 \sum_{p,j=1}^P \beta_p \beta_j \int_{(\mathbb{R}^d)^2} x_p x_j d\pi_{pj}(x_p, x_j), \quad (28)$$

where π_{pj} corresponds to marginalizing π onto all components but p, j . This can be formalized by defining the map $L_{pj}(x_1, \dots, x_P) = (x_p, x_j)$ and $\pi_{pj} = L_{pj\#}\pi$.

Now define $\pi^* = (C_{\mu_1}^{-1}, \dots, C_{\mu_P}^{-1})_{\#} \mathcal{U}(0, 1)$

Claim: π^* is optimal

First observe $L_{pj\#}\pi^* = (C_{\mu_p}^{-1}, C_{\mu_j}^{-1})_{\#} \mathcal{U}(0, 1)$ by marginalization. Note this is the optimal coupling between pairs μ_p, μ_j , see [25] (this can easily be obtained by observing that plugging in $(C_{\mu_p}^{-1}, C_{\mu_j}^{-1})_{\#} \mathcal{U}(0, 1)$ into the Wasserstein objective achieves the minimum – it is also a valid coupling, thus it has to be the optimal coupling.)

Now, note that

$$\arg \max_{\gamma \in \Pi(\mu_p, \mu_j)} \int_{(\mathbb{R}^d)^2} x_p x_j d\gamma = \arg \min_{\gamma \in \Pi(\mu_p, \mu_j)} \int_{(\mathbb{R}^d)^2} |x_p - x_j|^2 d\gamma, \quad (29)$$

and also that for any multimarginal coupling $\pi \in \Pi(\mu_1, \dots, \mu_P)$, π_{pj} is a pairwise coupling in $\Pi(\mu_p, \mu_j)$ by the transfer lemma.

We can hence deduce that $\forall \pi \in \Pi(\mu_1, \dots, \mu_P)$

$$\int_{(\mathbb{R}^d)^2} x_p x_j d\pi_{pj} \leq \int_{(\mathbb{R}^d)^2} x_p x_j d\pi_{pj}^* \quad \forall p, j = 1, \dots, P, \quad (30)$$

because both π_{pj} and π_{pj}^* are couplings of μ_p, μ_j and π_{pj}^* is optimal.

Therefore, it holds that

$$\int_{(\mathbb{R}^d)^P} \sum_{p=1}^P \beta_p \|x_p - \sum_j \beta_j x_j\|^2 d\pi^*(x_1, \dots, x_P) \quad (31)$$

$$= \sum_{p=1}^P \beta_p \int_{\mathbb{R}^d} |x_p|^2 d\mu_p - 2 \sum_{p,j=1}^P \beta_p \beta_j \int_{(\mathbb{R}^d)^2} x_p x_j d\pi_{pj}^*(x_p, x_j) \quad (32)$$

$$\leq \sum_{p=1}^P \beta_p \int_{\mathbb{R}^d} |x_p|^2 d\mu_p - 2 \sum_{p,j=1}^P \beta_p \beta_j \int_{(\mathbb{R}^d)^2} x_p x_j d\pi_{pj}(x_p, x_j) \quad (33)$$

$$= \int_{(\mathbb{R}^d)^P} \sum_{p=1}^P \beta_p \|x_p - \sum_j \beta_j x_j\|^2 d\pi(x_1, \dots, x_P), \quad (34)$$

which proves the claim that π^* is the optimal multi-marginal coupling. We now compute the distance by plugging in the optimal coupling:

$$\mathcal{MW}^2(\mu_1, \dots, \mu_P) = \int_{(\mathbb{R}^d)^P} \sum_{p=1}^P \beta_p |x_p - \sum_j \beta_j x_j|^2 d\pi^*(x_1, \dots, x_P) \quad (35)$$

$$= \int_{(\mathbb{R}^d)^P} \sum_{p=1}^P \beta_p |x_p - \sum_j \beta_j x_j|^2 d(C_{\mu_1}^{-1}, \dots, C_{\mu_P}^{-1})_{\#} \mathcal{U}(0, 1) \quad (36)$$

$$= \int_0^1 \sum_{p=1}^P \beta_p |C_{\mu_p}^{-1}(x) - \sum_j \beta_j C_{\mu_j}^{-1}(x)|^2 dx. \quad (37)$$

□

A.2 Generalized Metric Properties

Definition 9. Assume $\mu_p \in \mathcal{M}(\mathbb{R}^d)$, where $p = 1, \dots, P$, and let $D : \mathcal{M}(\mathbb{R}^d) \times \dots \times \mathcal{M}(\mathbb{R}^d) \rightarrow \mathbb{R}$ be a multi-marginal Wasserstein metric with barycentric weights β . Then, D is a generalized metric if the following properties hold:

1. $D(\mu_1, \dots, \mu_P) \geq 0$
2. $D(\mu_1, \dots, \mu_P) = 0 \Leftrightarrow \mu_1 = \dots = \mu_P$
3. $D(\mu_1, \dots, \mu_P) = D_{\sigma}(\mu_{\sigma(1)}, \dots, \mu_{\sigma(P)})$, $\forall \sigma \in \mathbb{S}_P$ where D_{σ} denotes that the barycentric weights β are permuted by σ and \mathbb{S}_P is the group of permutations of order P .
4. $\forall \mu \in \mathcal{M}(\mathbb{R}^d) : D(\mu_1, \dots, \mu_P) \leq \sum_{p=1}^P D(\mu_1, \dots, \mu_{p-1}, \mu, \mu_{p+1}, \dots, \mu_P)$

Proposition 10. \mathcal{MW} is a generalized metric on the restriction $\mathcal{M}(\mathbb{R})$.

Proof. Property (1), i.e., positivity is clear because

$$\mathcal{MW}^2(\mu_1, \dots, \mu_P) = \int_0^1 \sum_{p=1}^P \beta_p \left| C_{\mu_p}^{-1}(x) - \sum_{j=1}^P \beta_j C_{\mu_j}^{-1}(x) \right|^2 dx \geq 0 \quad (38)$$

Next, we prove property (2).

We begin by proving the forward implication (\Rightarrow).

$$\mathcal{MW}(\mu_1, \dots, \mu_P) = 0 \quad (39)$$

$$\Rightarrow \left(\int_0^1 \sum_{p=1}^P \beta_p \left| C_{\mu_p}^{-1}(x) - \sum_{j=1}^P \beta_j C_{\mu_j}^{-1}(x) \right|^2 dx \right)^{\frac{1}{2}} = 0 \quad (40)$$

$$\Rightarrow \int_0^1 \sum_{p=1}^P \beta_p \left| C_{\mu_p}^{-1}(x) - \sum_{j=1}^P \beta_j C_{\mu_j}^{-1}(x) \right|^2 dx = 0 \quad (41)$$

$$\Rightarrow C_{\mu_p}^{-1}(x) - \sum_{j=1}^P \beta_j C_{\mu_j}^{-1}(x) = 0 \quad \forall p = 1, \dots, P, \forall x \in [0, 1] \quad (42)$$

Now assume for contradiction that $\exists m, n, x : C_{\mu_m}^{-1}(x) \neq C_{\mu_n}^{-1}(x)$, then:

$$C_{\mu_m}^{-1}(x) = \sum_{j=1}^P \beta_j C_{\mu_j}^{-1}(x), \quad C_{\mu_n}^{-1}(x) = \sum_{j=1}^P \beta_j C_{\mu_j}^{-1}(x) \quad (43)$$

$$\Leftrightarrow C_{\mu_m}^{-1}(x) - C_{\mu_n}^{-1}(x) = \sum_{j=1}^P \beta_j C_{\mu_j}^{-1}(x) - \sum_{j=1}^P \beta_j C_{\mu_j}^{-1}(x) = 0 \quad (44)$$

which is a contradiction, therefore $C_{\mu_m}^{-1}(x) = C_{\mu_n}^{-1}(x) \quad \forall m, n, x$, thus $\mu_1 = \dots = \mu_P$

We continue by proving the backward implication (\Leftarrow).

If $\mu_1 = \dots = \mu_P$, then $C_{\mu_p}^{-1}(x) = C_{\mu_{p'}}^{-1}(x) \quad \forall x, \forall p, p' = 1, \dots, P$.

Therefore, $C_{\mu_p}^{-1}(x) - \sum_{j=1}^P \beta_j C_{\mu_j}^{-1}(x) = 0 \quad \forall p = 1, \dots, P, \forall x \in [0, 1]$. Thus,

$$\mathcal{MW}(\mu_1, \dots, \mu_P) = \left(\int_0^1 \sum_{p=1}^P \beta_p \left| C_{\mu_p}^{-1}(x) - \sum_{j=1}^P \beta_j C_{\mu_j}^{-1}(x) \right|^2 dx \right)^{\frac{1}{2}} = 0. \quad (45)$$

We continue with permutation invariance (3),

$$\mathcal{MW}(\mu_1, \dots, \mu_P) = \left(\int_0^1 \sum_{p=1}^P \beta_p \left| C_{\mu_p}^{-1}(x) - \sum_{j=1}^P \beta_j C_{\mu_j}^{-1}(x) \right|^2 dx \right)^{\frac{1}{2}} \quad (46)$$

$$= \left(\int_0^1 \sum_{p=1}^P \beta_p \left| C_{\mu_p}^{-1}(x) - \sum_{j=1}^P \beta_{\sigma(j)} C_{\mu_{\sigma(j)}}^{-1}(x) \right|^2 dx \right)^{\frac{1}{2}} \quad (47)$$

$$= \left(\int_0^1 \sum_{p=1}^P \beta_{\sigma(p)} \left| C_{\mu_{\sigma(p)}}^{-1}(x) - \sum_{j=1}^P \beta_j C_{\mu_{\sigma(j)}}^{-1}(x) \right|^2 dx \right)^{\frac{1}{2}} \quad (48)$$

$$= \mathcal{MW}_{\sigma}(\mu_{\sigma(1)}, \dots, \mu_{\sigma(P)}) \quad (49)$$

Equalities holds because sums are invariant under any permutation σ .

We finally prove the generalized triangle inequality (4). Note the slight abuse of notation that $p + 1$ component does not exist when $p = P$.

We begin by proving the case $P \geq 3$. Firstly, we rewrite the multi-marginal functional in the following way:

$$\mathcal{MW}^2(\mu_1, \dots, \mu_P) = \sum_{p=1}^P \beta_p \int_0^1 \left| C_{\mu_p}^{-1}(x) - \sum_{j=1}^P \beta_j C_{\mu_j}^{-1}(x) \right|^2 dx \quad (50)$$

$$= \frac{1}{2} \sum_{p, p'=1}^P \beta_p \beta_{p'} \int_0^1 \left| C_{\mu_p}^{-1}(x) - C_{\mu_{p'}}^{-1}(x) \right|^2 dx \quad (51)$$

$$= \frac{1}{2} \sum_{p, p'=1}^P \beta_p \beta_{p'} \int_0^1 f_{p, p'}^2(x) dx \quad (52)$$

where $f_{p,p'}(x) = \left| C_{\mu_p}^{-1}(x) - C_{\mu_{p'}}^{-1}(x) \right|$. The results holds because

$$\sum_{m,n=1}^P \beta_m \beta_n |C_{\mu_m}^{-1}(x) - C_{\mu_n}^{-1}(x)|^2 = \sum_{m=1}^P \beta_m \left| C_{\mu_m}^{-1}(x) - \sum_{n=1}^P \beta_n C_{\mu_j}^{-1}(x) \right|^2, \quad (53)$$

which holds because

$$\sum_{m=1}^P \beta_m \left| x_m - \sum_{n=1}^P \beta_n x_n \right|^2 \quad (54)$$

$$= \sum_{m=1}^P \beta_m \left[|x_m|^2 + \left| \sum_{n=1}^P \beta_n x_n \right|^2 - 2 \sum_{n=1}^P \beta_n x_m x_n \right] \quad (55)$$

$$= \sum_{m=1}^P \beta_m |x_m|^2 + \sum_{m,n=1}^P \beta_m \beta_n x_m x_n - 2 \sum_{m,n=1}^P \beta_m \beta_n x_m x_n \quad (56)$$

$$= \sum_{m=1}^P \beta_m |x_m|^2 - \sum_{m,n=1}^P \beta_m \beta_n x_m x_n \quad (57)$$

$$= \sum_{m,n=1}^P \beta_m \beta_n |x_m|^2 - \sum_{m,n=1}^P \beta_m \beta_n x_m x_n \quad (58)$$

$$= \sum_{m,n=1}^P \beta_m \beta_n \left(\frac{1}{2} |x_m|^2 + \frac{1}{2} |x_n|^2 - x_m x_n \right) \quad (59)$$

$$= \frac{1}{2} \sum_{m,n=1}^P \beta_m \beta_n |x_m - x_n|^2. \quad (60)$$

Therefore, we have

$$\sum_{p=1}^P \mathcal{M}\mathcal{W}^2(\mu_1, \dots, \mu_{p-1}, \mu, \mu_{p+1}, \dots, \mu_P) = \frac{1}{2} \sum_{p=1}^P \sum_{m,n \neq p}^P \beta_m \beta_n \int_0^1 f_{n,m}^2(x) dx + C, \quad (61)$$

where $C > 0$.

We now show that $\int_0^1 \sum_{p=1}^P \sum_{m,n \neq p}^P \beta_m \beta_n f_{n,m}^2(x) dx \geq \sum_{p,p'=1}^P \beta_p \beta_{p'} \int_0^1 f_{p,p'}^2(x) dx$. This can be observed by noting that all $\int_0^1 f_{p,p'}^2(x) dx$ terms on the RHS appear on the LHS. Indeed, for any m', n' , $\int_0^1 f_{m',n'}^2(x) dx$ appears in the $p \neq m', n'$ summation, which always holds for some p as $P \geq 3$.

Therefore, we have shown that

$$\mathcal{M}\mathcal{W}^2(\mu_1, \dots, \mu_P) \leq \sum_{p=1}^P \mathcal{M}\mathcal{W}^2(\mu_1, \dots, \mu_{p-1}, \mu, \mu_{p+1}, \dots, \mu_P) \quad (62)$$

Also,

$$\mathcal{M}\mathcal{W}^2(\mu_1, \dots, \mu_P) \leq \sum_{p=1}^P \mathcal{M}\mathcal{W}^2(\mu_1, \dots, \mu_{p-1}, \mu, \mu_{p+1}, \dots, \mu_P) \quad (63)$$

$$\Rightarrow \mathcal{M}\mathcal{W}(\mu_1, \dots, \mu_P) \leq \sqrt{\sum_{p=1}^P \mathcal{M}\mathcal{W}^2(\mu_1, \dots, \mu_{p-1}, \mu, \mu_{p+1}, \dots, \mu_P)} \quad (64)$$

$$\leq \sum_{p=1}^P \sqrt{\mathcal{M}\mathcal{W}^2(\mu_1, \dots, \mu_{p-1}, \mu, \mu_{p+1}, \dots, \mu_P)} \quad (65)$$

$$= \sum_{p=1}^P \mathcal{M}\mathcal{W}(\mu_1, \dots, \mu_{p-1}, \mu, \mu_{p+1}, \dots, \mu_P) \quad (66)$$

which proves the result. The case $P = 2$ has been proved via different approaches (e.g. [25]).

□

Proposition 4. \mathcal{SMW} is a generalized metric on the restriction $\mathcal{M}(\mathbb{R}^d)$.

Proof. Property (1) holds by definition due to positivity of \mathcal{MW} on \mathbb{R} and the definition of the sliced multi-marginal distance.

Property (2) is more delicate. We begin with the forward direction (\Rightarrow).

We extend the proof of Nadjahi et al. [21] to the multi-marginal case. Define Θ as the uniform distribution on S_{d-1} . Define ‘for (Θ -almost-every) θ ’ as $\forall \Theta$ -a-e- θ . Firstly, the following holds:

$$\mathcal{SMW}(\mu_1, \dots, \mu_P) = 0 \quad (67)$$

$$\Rightarrow \left(\frac{1}{\text{Vol}(S_{d-1})} \int_{S_{d-1}} \mathcal{MW}^2(M_{\theta\#}\mu_1, \dots, M_{\theta\#}\mu_P) d\Theta(\theta) \right)^{\frac{1}{2}} = 0 \quad (68)$$

$$\Rightarrow \mathcal{MW}(M_{\theta\#}\mu_1, \dots, M_{\theta\#}\mu_P) \forall \Theta\text{-a-e-}\theta \quad (69)$$

$$\Rightarrow M_{\theta\#}\mu_1 = \dots = M_{\theta\#}\mu_P \forall \Theta\text{-a-e-}\theta \quad (70)$$

Next, we define the Fourier transform of any measure μ on $\mathcal{M}(\mathbb{R}^s)$, $s \geq 1$ at any $\mathbf{w} \in \mathbb{R}^s$:

$$\mathcal{F}[\mu](\mathbf{w}) = \int_{\mathbb{R}^s} e^{-i\langle \mathbf{w}, \mathbf{x} \rangle} d\mu(\mathbf{x}). \quad (71)$$

Therefore, using properties of push-forwards, the following holds:

$$\mathcal{F}[M_{\theta\#}\mu](t) = \int_{\mathbb{R}} e^{-itu} dM_{\theta\#}\mu(u) = \int_{\mathbb{R}^s} e^{-it\langle \theta, \mathbf{x} \rangle} d\mu(\mathbf{x}) = \mathcal{F}[\mu](t\theta). \quad (72)$$

As $\forall \Theta$ -a-e- θ , $M_{\theta\#}\mu_1 = \dots = M_{\theta\#}\mu_P$, then $\mathcal{F}[M_{\theta\#}\mu_1] = \dots = \mathcal{F}[M_{\theta\#}\mu_P]$, which implies that $\mathcal{F}[\mu_1] = \dots = \mathcal{F}[\mu_P]$. By injectivity of the Fourier transform, we conclude that $\mu_1 = \dots = \mu_P$.

We continue with the backward direction (\Leftarrow).

We assume $\mu_1 = \dots = \mu_P$, which implies the following:

$$\mu_1 = \dots = \mu_P \quad (73)$$

$$\Rightarrow M_{\theta\#}\mu_1 = \dots = M_{\theta\#}\mu_P \forall \Theta\text{-a-e-}\theta \quad (74)$$

$$\Rightarrow \mathcal{MW}^2(M_{\theta\#}\mu_1, \dots, M_{\theta\#}\mu_P) = 0 \forall \Theta\text{-a-e-}\theta \quad (75)$$

$$\Rightarrow \mathcal{SMW}(\mu_1, \dots, \mu_P) = \left(\frac{1}{\text{Vol}(S_{d-1})} \int_{S_{d-1}} \mathcal{MW}^2(M_{\theta\#}\mu_1, \dots, M_{\theta\#}\mu_P) d\Theta(\theta) \right)^{\frac{1}{2}} = 0. \quad (76)$$

We now prove Property (3)

$$\mathcal{SMW}(\mu_1, \dots, \mu_P) = \left(\frac{1}{\text{Vol}(S_{d-1})} \int_{S_{d-1}} \mathcal{MW}^2(M_{\theta\#}\mu_1, \dots, M_{\theta\#}\mu_P) d\Theta(\theta) \right)^{\frac{1}{2}} \quad (77)$$

$$= \left(\frac{1}{\text{Vol}(S_{d-1})} \int_{S_{d-1}} \mathcal{MW}_\sigma^2(M_{\theta\#}\mu_{\sigma(1)}, \dots, M_{\theta\#}\mu_{\sigma(P)}) d\Theta(\theta) \right)^{\frac{1}{2}} \quad (78)$$

$$= \mathcal{SMW}_\sigma(\mu_{\sigma(1)}, \dots, \mu_{\sigma(P)}) \quad (79)$$

We finally end by proving Property (4), the generalized triangle inequality.

Earlier, we showed that

$$\mathcal{MW}^2(\mu_1, \dots, \mu_P) \leq \sum_{p=1}^P \mathcal{MW}^2(\mu_1, \dots, \mu_{p-1}, \mu, \mu_{p+1}, \dots, \mu_P). \quad (80)$$

This implies that

$$\mathcal{SMW}^2(\mu_1, \dots, \mu_P) \quad (81)$$

$$= \frac{1}{\text{Vol}(S_{d-1})} \int_{S_{d-1}} \mathcal{MW}^2(M_{\theta\#}\mu_1, \dots, M_{\theta\#}\mu_P) d\Theta(\theta) \quad (82)$$

$$\leq \sum_{p=1}^P \frac{1}{\text{Vol}(S_{d-1})} \int_{S_{d-1}} \mathcal{MW}^2(M_{\theta\#}\mu_1, \dots, M_{\theta\#}\mu_{p-1}, M_{\theta\#}\mu, M_{\theta\#}\mu_{p+1}, \dots, M_{\theta\#}\mu_P) d\Theta(\theta) \quad (83)$$

$$= \sum_{p=1}^P \mathcal{SMW}^2(\mu_1, \dots, \mu_{p-1}, \mu, \mu_{p+1}, \dots, \mu_P). \quad (84)$$

Therefore, we conclude that

$$\mathcal{SMW}^2(\mu_1, \dots, \mu_P) \leq \sum_{p=1}^P \mathcal{SMW}^2(\mu_1, \dots, \mu_{p-1}, \mu, \mu_{p+1}, \dots, \mu_P) \quad (85)$$

$$\Rightarrow \mathcal{SMW}(\mu_1, \dots, \mu_P) \leq \sum_{p=1}^P \mathcal{SMW}(\mu_1, \dots, \mu_{p-1}, \mu, \mu_{p+1}, \dots, \mu_P) \quad (86)$$

directly in the same way as in the proof of Proposition the generalized triangle inequality for \mathcal{MW} . \square

A.3 Mathematical Properties

Proposition 3.

$$\mathcal{SMW}^2(\mu_1, \dots, \mu_P) = \frac{1}{2} \sum_{i,j=1}^P \beta_i \beta_j \mathcal{SW}^2(\mu_i, \mu_j) \quad (87)$$

Proof.

$$\mathcal{SMW}^2(\mu_1, \dots, \mu_P) = \frac{1}{\text{Vol}(S_{d-1})} \int_{S_{d-1}} \int_{\mathbb{R}^d} \frac{1}{2} \sum_{i,j=1}^P \beta_i \beta_j |x_i - x_j|^2 d\pi^{*\theta}(x_1, \dots, x_P) d\Theta(\theta) \quad (88)$$

$$= \frac{1}{2\text{Vol}(S_{d-1})} \sum_{i,j=1}^P \beta_i \beta_j \int_{S_{d-1}} \int_{\mathbb{R} \times \mathbb{R}} |x_i - x_j|^2 d\pi_{ij}^{*\theta}(x_i, x_j) d\Theta(\theta) \quad (89)$$

$$= \frac{1}{2} \sum_{i,j=1}^P \beta_i \beta_j \frac{1}{\text{Vol}(S_{d-1})} \int_{S_{d-1}} \mathcal{W}^2(M_{\theta\#}\mu_i, M_{\theta\#}\mu_j) d\Theta(\theta), \quad (90)$$

where $\pi^{*\theta}$ is the optimal coupling between $M_{\theta\#}\mu_1, \dots, M_{\theta\#}\mu_P$ and $M_{\theta}(x) = \langle x, \theta \rangle$. Similarly to proofs of closed-form formulas for multi-marginal Kantorovich transport, we know that $\pi_{ij}^{*\theta}$ is the optimal coupling between $M_{\theta\#}\mu_i, M_{\theta\#}\mu_j$. As a result, it holds that

$$\mathcal{SMW}^2(\mu_1, \dots, \mu_P) = \frac{1}{2} \sum_{i,j=1}^P \beta_i \beta_j \mathcal{SW}^2(\mu_i, \mu_j). \quad (91)$$

\square

Corollary 11.

$$\mathcal{SMW}^2(\mu_1, \dots, \mu_P) \leq \frac{1}{2} \sum_{i,j=1}^P \beta_i \beta_j \mathcal{W}^2(\mu_i, \mu_j) \quad (92)$$

Proof. By Proposition 3, it holds that

$$\mathcal{SMW}^2(\mu_1, \dots, \mu_P) = \frac{1}{2} \sum_{i,j=1}^P \beta_i \beta_j \mathcal{SW}^2(\mu_i, \mu_j). \quad (93)$$

Also, Bonnotte [9] shows that

$$\mathcal{SW}^2(\mu, \nu) \leq \mathcal{W}^2(\mu, \nu) \quad \forall \mu, \nu. \quad (94)$$

The result follows directly. \square

Corollary 12.

$$0 \leq a^{2(d+1)} \sum_{i,j=1}^P \beta_i \beta_j \mathcal{W}^{4(d+1)}(\mu_i, \mu_j) \leq b^{2(d+1)} \mathcal{SMW}^2(\mu_1, \dots, \mu_P). \quad (95)$$

Proof. Bonnotte [9] has shown that it holds for some positive constants a, b that

$$0 \leq a \mathcal{W}^2(\mu_i, \mu_j) \leq b \mathcal{SW}^{1/(d+1)}(\mu_i, \mu_j), \quad (96)$$

and that x^{d+1} is an increasing function for all positive x . Therefore, raising both sides to the power of $2(d+1)$, we obtain that

$$0 \leq a^{2(d+1)} \mathcal{W}^{4(d+1)}(\mu_i, \mu_j) \leq b^{2(d+1)} \mathcal{SW}^2(\mu_i, \mu_j). \quad (97)$$

Now summing across i, j , and weighting with the barycentric cost's weights, we obtain

$$0 \leq a^{2(d+1)} \sum_{i,j=1}^P \beta_i \beta_j \mathcal{W}^{4(d+1)}(\mu_i, \mu_j) \leq b^{2(d+1)} \mathcal{SMW}^2(\mu_1, \dots, \mu_P). \quad (98)$$

It therefore follows that as $\mathcal{SMW}^2(\mu_1, \dots, \mu_P) \rightarrow 0$, we also have $\mathcal{W}^{4(d+1)}(\mu_i, \mu_j) \rightarrow 0$ for each pair of measures, and hence by positivity of \mathcal{W} that $\sum_{i,j=1}^P \beta_i \beta_j \mathcal{W}^2(\mu_i, \mu_j) \rightarrow 0$. \square

A.4 Sample/Projection Complexity

We now study $E[\mathcal{SMW}^2(\mu_1, \dots, \mu_P) - \mathcal{SMW}^2(\hat{\mu}_1, \dots, \hat{\mu}_P)]^2$ where $\hat{\mu}_p$'s refers to empirical measures with n samples. Then the following result holds:

Proposition 5. *If $\mu_1, \dots, \mu_P \in \mathcal{M}(\mathbb{R}^d)$, and assuming \mathcal{W}^2 has sample complexity $\rho(N)$ on \mathbb{R} , then*

$$E[\mathcal{SMW}^2(\mu_1, \dots, \mu_P) - \mathcal{SMW}^2(\hat{\mu}_1, \dots, \hat{\mu}_P)]^2 \leq \frac{1}{2} \rho(N). \quad (99)$$

This result shows the sample complexity is dimension free.

Proof. We conclude from Proposition 3

$$\mathcal{SMW}^2(\mu_1, \dots, \mu_P) - \mathcal{SMW}^2(\hat{\mu}_1, \dots, \hat{\mu}_P) = \frac{1}{2} \sum_{i,j=1}^P \beta_i \beta_j \left(\mathcal{SW}^2(\mu_i, \mu_j) - \mathcal{SW}^2(\hat{\mu}_i, \hat{\mu}_j) \right). \quad (100)$$

If \mathcal{W}^2 on \mathbb{R} has sample complexity $\rho(N)$, then \mathcal{SW}^2 on \mathbb{R}^d also has sample complexity $\rho(N)$, i.e., its sample complexity is dimension free. The proof relies on an application of Jensen's inequality and is a special case of Nadjahi et al. [22].

$$E \left| \mathcal{SW}^2(\mu, \nu) - \mathcal{SW}^2(\hat{\mu}_n, \hat{\nu}_n) \right| = E \left| \int_{S_{d-1}} \{ \mathcal{W}^2(\theta_{\#}^* \mu, \theta_{\#}^* \nu) - \mathcal{W}^2(\theta_{\#}^* \hat{\mu}_n, \theta_{\#}^* \hat{\nu}_n) \} d\Theta(\theta) \right| \quad (101)$$

$$\leq E \left\{ \int_{S_{d-1}} \left| \mathcal{W}^2(\theta_{\#}^* \mu, \theta_{\#}^* \nu) - \mathcal{W}^2(\theta_{\#}^* \hat{\mu}_n, \theta_{\#}^* \hat{\nu}_n) \right| d\Theta(\theta) \right\} \quad (102)$$

$$\leq \int_{S_{d-1}} E \left| \mathcal{W}^2(\theta_{\#}^* \mu, \theta_{\#}^* \nu) - \mathcal{W}^2(\theta_{\#}^* \hat{\mu}_n, \theta_{\#}^* \hat{\nu}_n) \right| d\Theta(\theta) \quad (103)$$

$$\leq \int_{S_{d-1}} \rho(N) d\Theta(\theta) = \rho(N) \quad (104)$$

Hence,

$$E \left| \mathcal{SMW}^2(\mu_1, \dots, \mu_P) - \mathcal{SMW}^2(\hat{\mu}_1, \dots, \hat{\mu}_P) \right| \quad (105)$$

$$= E \left| \frac{1}{2} \sum_{i,j=1}^P \beta_i \beta_j \left(\mathcal{SW}^2(\mu_i, \mu_j) - \mathcal{SW}^2(\hat{\mu}_i, \hat{\mu}_j) \right) \right| \quad (106)$$

$$\leq \frac{1}{2} \sum_{i,j=1}^P \beta_i \beta_j E \left| \mathcal{SW}^2(\mu_i, \mu_j) - \mathcal{SW}^2(\hat{\mu}_i, \hat{\mu}_j) \right| \quad (107)$$

$$\leq \frac{1}{2} \sum_{i,j=1}^P \beta_i \beta_j \rho(N) = \frac{1}{2} \rho(N). \quad (108)$$

□

Here we also derive similar results to theirs about projection complexity.

Proposition 6. *Let $\mu_1, \dots, \mu_P \in \mathcal{M}(\mathbb{R}^d)$, and define $\overline{\mathcal{SMW}}$ the approximation obtained by uniformly picking L projections on S_{d-1} , then*

$$\mathbb{E} \left[\overline{\mathcal{SMW}}^2(\mu_1, \dots, \mu_P) - \mathcal{SMW}^2(\mu_1, \dots, \mu_P) \right]^2 \leq L^{-1/2} \text{Var}_{\theta} \left[\mathcal{MW}^2(\mu_1^{\theta}, \dots, \mu_P^{\theta}) \right], \quad (109)$$

where θ follows the uniform distribution on S_{d-1} and $\mu_p^{\theta} = M_{\#}^{\theta}(\mu_p)$.

Proof. We bound the error arising from the Monte Carlo approximation of \mathcal{SMW} , similarly to Nadjahi et al. [22] in the pairwise case. In particular, define $\delta = \int_{S_{d-1}} \mathcal{MW}^2(M_{\theta\#}\mu_1, \dots, M_{\theta\#}\mu_P) d\Theta(\theta)$. Then we have that

$$E_{\theta \sim \sigma} \left| \overline{\mathcal{SMW}}^2(\mu_1, \dots, \mu_P) - \mathcal{SMW}^2(\mu_1, \dots, \mu_P) \right| \quad (110)$$

$$\leq \left\{ E_{\theta \sim \sigma} \left| \overline{\mathcal{SMW}}^2(\mu_1, \dots, \mu_P) - \mathcal{SMW}^2(\mu_1, \dots, \mu_P) \right|^2 \right\}^{\frac{1}{2}} \quad (111)$$

$$\leq L^{-1/2} \int_{S_{d-1}} \left\{ \mathcal{MW}^2(M_{\theta\#}\mu_1, \dots, M_{\theta\#}\mu_P) - \delta \right\}^2 d\Theta(\theta) \quad (112)$$

$$= L^{-1/2} \text{Var}_{\theta} \left[\mathcal{MW}^2(\mu_1^{\theta}, \dots, \mu_P^{\theta}) \right], \quad (113)$$

which holds due to the same Monte-Carlo concentration inequality as in Nadjahi et al. [22] (Proof of Theorem 6). □

A.5 Equivalence to Sliced Barycenters and Weak Convergence

Proposition 7. *Let $\mu_1, \dots, \mu_P \in \mathcal{M}(\mathbb{R}^d)$, $\sum_{p=1}^P \beta_p = 1$. Furthermore, let $\hat{\beta}_p$ be augmented multi-marginal weights, so that for $m \in [0, 1]$ it holds that $\hat{\beta}_p = m\beta_p$ for $p = 1, \dots, P$, $\sum_{p=1}^{P+1} \hat{\beta}_p = 1$,*

and $\mathcal{D} = \mathcal{SW}^2$. Then

$$\arg \min_{\mu \in \mathcal{M}(\mathbb{R}^d)} \mathcal{SMW}^2(\mu_1, \dots, \mu_P, \mu) = \arg \min_{\mu \in \mathcal{M}(\mathbb{R}^d)} \mathcal{F}(\mu), \quad (114)$$

where β is the weight vector of \mathcal{F} and $\hat{\beta}$ is the weight vector of \mathcal{SMW} .

Proof.

$$\arg \min_{\mu \in \mathcal{M}(\mathbb{R}^d)} \mathcal{SMW}^2(\mu_1, \dots, \mu_P, \mu) \quad (115)$$

$$= \arg \min_{\mu \in \mathcal{M}(\mathbb{R}^d)} \sum_{p=1}^P \hat{\beta}_p \hat{\beta}_{P+1} \mathcal{SW}^2(\mu, \mu_p) \quad (116)$$

$$= \arg \min_{\mu \in \mathcal{M}(\mathbb{R}^d)} \sum_{p=1}^P \beta_p \mathcal{SW}^2(\mu_p, \mu) \quad (117)$$

$$= \arg \min_{\mu \in \mathcal{M}(\mathbb{R}^d)} \mathcal{F}(\mu). \quad (118)$$

□

A.6 Differentiability

Proposition 8. Let $\mu_1, \dots, \mu_P \in \mathcal{M}(\mathbb{R}^d)$ be discrete measures with N atoms, which we gather into matrices $\{\mathbf{X}^{(p)}\}_{p=1}^P$, and similarly define $\mu_{\mathbf{X}}$ with atoms \mathbf{X} . Assume \mathbf{X} has distinct points. Then \mathcal{SMW}^2 is smooth with gradient

$$\nabla_{\mathbf{X}} \mathcal{SMW}^2(\mu_1, \dots, \mu_P, \mu_{\mathbf{X}}) = \beta_{P+1} \sum_{p=1}^P \beta_p \int_{S_{d-1}} \mathbf{X}_{\theta} - (\mathbf{X}_{\theta}^{(p)} \circ \sigma_{\mathbf{X}_{\theta}} \circ \sigma_{\mathbf{X}_{\theta}^{(p)}}^{-1}) d\theta, \quad (119)$$

where $\sigma_{\mathbf{X}}$ is the permutation that sorts atoms of \mathbf{X} .

Proof. Define $\sigma_{\mathbf{Y}}$ be the permutation of $\{1, \dots, N\}$ that sorts atoms of \mathbf{Y} . Also, define $\mathbf{X}_{\theta} \in \mathbb{R}^N$, such that $(\mathbf{X}_{\theta})_i = \langle x_i, \theta \rangle$. Then

$$\mathcal{SMW}^2(\mu_1, \dots, \mu_P, \mu_{\mathbf{X}}) = \sum_{p=1}^P \beta_{P+1} \beta_p \mathcal{SW}^2(\mu_{\mathbf{X}}, \mu_p) + C(\mu_1, \dots, \mu_P). \quad (120)$$

Hence,

$$\nabla_{\mathbf{X}} \mathcal{SMW}^2(\mu_1, \dots, \mu_P, \mu_{\mathbf{X}}) = \nabla_{\mathbf{X}} \sum_{p=1}^P \beta_{P+1} \beta_p \mathcal{SW}^2(\mu_{\mathbf{X}}, \mu_p) \quad (121)$$

$$= \sum_{p=1}^P \beta_{P+1} \beta_p \int_{S_{d-1}} \mathbf{X}_{\theta} - \mathbf{X}_{\theta}^{(p)} \circ (\sigma_{\mathbf{X}_{\theta}} \circ \sigma_{\mathbf{X}_{\theta}^{(p)}}^{-1}) d\theta. \quad (122)$$

The last equality is due to Bonneel et al. [8].

□

B Additional Experimental Details

We now provide further experimental details. All experiments ran on CPU, besides the benchmarking experiments, which ran on a single P100 GPU.

Ellipses - Multi-Task Density Estimation

We set the batch size to 150, and parametrize each measure ν_p as a discrete measure with 150 atoms which we optimize over via stochastic gradient descent. We set the number of projections to 20.

Multi-Task Reinforcement Learning

The horizon is set to $T = 200$. The learning rate is set to 2.5×10^{-4} , and the batch size to optimize the Q -function to 32. The Q-network is a 2-layer MLP with tanh activation. We use $f(x) = e^{-5x}$ to rescale the reward function following Dadashi et al. [15], we set the number of projections to $K = 50$ and $\gamma = 1$. Also, we set $\alpha = \frac{1}{30}$. Our implementation extends the repository <https://github.com/xtma/simple-pytorch-rl> to the multi-task setting, and leverages OpenAI gym environments [10].

Gradient Flow experiment

We follow the setup of Bonneel et al. [8]. In particular, we discretize the flow to numerically estimate it via gradient descent $\mathbf{X}^{(l+1)} = \mathbf{X}^{(l)} - \nabla \mathcal{SMW}^2(\mu_1, \dots, \mu_P, \mu_{\mathbf{X}^{(l)}})$, and plot the location of particles for $l = 0, \dots, T$ where T is the number of steps (200), which approximates the gradient flow. We estimate \mathcal{SMW} with 30 projections. Each measure (including the initial measure μ_0) consist in samples from isotropic Gaussians, and the initial measure.