

Slicing Unbalanced Optimal Transport

Clément Bonet*

CREST, ENSAE, IP Paris, Palaiseau, France

clement.bonet@ensae.fr

Kimia Nadjahi*

CNRS, ENS, Paris, France

kimia.nadjahi@ens.fr

Thibault Séjourné*

LTS4, EPFL, Lausanne, Switzerland

thibault.sejourne@epfl.ch

Kilian Fatras

Mila, McGill University, Montreal, Canada

kilian.fatras@mila.quebec

Nicolas Courty

IRISA, Université Bretagne-Sud, Vannes, France

nicolas.courty@irisa.fr

Reviewed on OpenReview: <https://openreview.net/forum?id=AjJTg5MOr8>

Abstract

Optimal transport (OT) is a powerful framework to compare probability measures, a fundamental task in many statistical and machine learning problems. Substantial advances have been made in designing OT variants which are either computationally and statistically more efficient or robust. Among them, sliced OT distances have been extensively used to mitigate optimal transport’s cubic algorithmic complexity and curse of dimensionality. In parallel, unbalanced OT was designed to allow comparisons of more general positive measures, while being more robust to outliers. In this paper, we bridge the gap between those two concepts and develop a general framework for efficiently comparing positive measures. We notably formulate two different versions of sliced unbalanced OT, and study the associated topology and statistical properties. We then develop a GPU-friendly Frank-Wolfe like algorithm to compute the corresponding loss functions, and show that the resulting methodology is modular as it encompasses and extends prior related work. We finally conduct an empirical analysis of our loss functions and methodology on both synthetic and real datasets, to illustrate their computational efficiency, relevance and applicability to real-world scenarios including geophysical data.

1 Introduction

Many machine learning tasks involve aligning objects such as images, graphs, datasets or their representations after transformations. This is particularly relevant in transfer learning tasks like domain adaptation (Fatras et al., 2021) or multimodal machine learning (Baltrušaitis et al., 2018). These objects can be conveniently represented as positive measures, *i.e.*, a set of samples associated with non-negative weights. Aligning then consists in minimizing a distance or discrepancy between two measures. It is crucial to choose a meaningful discrepancy that has desirable statistical, robustness and computational properties. In particular, some settings require comparing arbitrary positive measures, *i.e.*, measures whose total mass can have an arbitrary value, as opposed to probability distributions whose total mass is equal to 1. In cell biology, for instance, measures are used to represent and compare gene expressions of cell populations, and the total mass corresponds to the population size (Schiebinger et al., 2019).

* Equal contribution

(Unbalanced) Optimal Transport. Optimal transport (OT) has been frequently chosen as a loss function to align objects. OT defines a distance between two positive measures α and β of the same mass ($m(\alpha) = m(\beta)$) by moving the mass of α toward the mass of β with the least possible effort. However, in some applications, the mass equality constraint is not satisfied, *i.e.*, $m(\alpha) \neq m(\beta)$. It can still be enforced by a re-normalization of the mass, which is potentially spurious and makes the problem less interpretable. This setting has motivated the development of a new OT framework, called *unbalanced OT* (UOT), that can naturally compare measures of different masses by softly relaxing the mass conservation constraints (Kondratyev et al., 2016; Liero et al., 2018; Chizat et al., 2018b). An appealing outcome of this new OT variant is its robustness to outliers which is achieved by discarding them before transporting α to β . UOT has been useful for many theoretical and practical applications, *e.g.*, theory of deep learning (Chizat & Bach, 2018; Rotskoff et al., 2019), biology (Schiebinger et al., 2019; Demetci et al., 2022) and domain adaptation (Fatras et al., 2021). We refer to (Séjourné et al., 2022a) for an extensive survey of UOT. Computing UOT requires to solve a linear program whose complexity is cubical in the number n of samples ($\mathcal{O}(n^3 \log n)$) (Pele & Werman, 2009; Peyré et al., 2019). Besides, accurately estimating UOT distances through empirical distributions is challenging as they suffer from the curse of dimension (Dudley, 1969). A common workaround is to rely on variants with lower complexities and better statistical properties. Among the most popular, we can list entropic OT (Cuturi, 2013; Pham et al., 2020) or minibatch OT (Fatras et al., 2020; 2021). In this paper, we focus on developing sliced UOT approaches.

Sliced Optimal Transport. Sliced OT (SOT) defines an alternative metric by leveraging the closed-form solution of OT between univariate measures (Rabin et al., 2012; Bonneel et al., 2015). It averages the OT cost between projections of (α, β) on 1D subspaces of \mathbb{R}^d . For 1D data, the OT solution can be computed through a sort algorithm, leading to an appealing $\mathcal{O}(n \log(n))$ complexity (Peyré et al., 2019). Furthermore, it has been shown to lift useful topological and statistical properties of OT from 1-dimensional to multi-dimensional settings (Nadjahi et al., 2020b; Bayraktar & Guo, 2021; Goldfeld & Greenewald, 2021). It therefore helps to mitigate the curse of dimensionality making SOT-based algorithms theoretically grounded, statistically efficient, and practical to solve even on large-scale settings. These appealing properties motivated the development of several variants and generalizations, *e.g.*, by considering different types or distributions of projections (Kolouri et al., 2019; Deshpande et al., 2019; Nguyen et al., 2020; Ohana et al., 2023; Nguyen et al., 2023) or manifold data (Bonet et al., 2023a;b;c). Fast computations of partial OT (a particular case of UOT) between univariate measures (Bonneel & Coeurjolly, 2019; Bai et al., 2023) or more generally on trees (Sato et al., 2020; Le & Nguyen, 2021), have been developed, so that slicing partial OT benefits from these efficient implementations and allows to compare large unnormalized measures. However, while (sliced) partial OT allows to compare measures with different masses, it assumes that each input measure is discrete and supported on points that all share the same mass (typically 1). In contrast, the Gaussian-Hellinger-Kantorovich (GHK) distance (Liero et al., 2018), another popular formulation of UOT, allows to compare measures with different masses *and* supported on points with varying masses, and has not been studied jointly with slicing.

Contributions. In this paper, we present the first general framework combining UOT and slicing between arbitrary distributions. Our main contribution is the introduction of two novel sliced variants of UOT, called *Sliced UOT* (SUOT) and *Unbalanced Sliced OT* (USOT). SUOT and USOT are both defined as regularized OT problems which leverage one-dimensional projections, but differ on how they relax the mass preservation constraint: USOT essentially performs a global reweighting of the inputs measures (α, β) , while SUOT reweights each projection of (α, β) . We provide a theoretical analysis of SUOT and USOT, which reveals that they share topological properties with UOT while being statistically more efficient in high-dimensional regimes, thanks to the slicing operation. Additionally, we propose fast and GPU-friendly algorithms to compute SUOT and USOT, based on the (non-trivial) dual derivation of our SUOT and USOT losses and a Frank-Wolfe strategy (Séjourné et al., 2022b). Finally, we illustrate the efficiency of our framework on various experiments: we deploy SUOT and USOT to compare distributions on non-Euclidean hyperbolic manifolds, classify documents, transfer image colors and aggregate large-scale geophysical data, and discuss their advantages over existing approaches.

Outline. In Section 2, we provide background knowledge on unbalanced OT and sliced OT. In Section 3, we introduce our new loss functions (SUOT and USOT) and prove their metric, topological, statistical and

duality properties in wide generality. We then explain in Section 4 how to compute SUOT and USOT via a Frank-Wolfe-based approach. We finally analyze the performance of SUOT and USOT on different practical tasks in Section 5.

2 Background

In this section, we first state our notations. Then, we provide the necessary background on unbalanced optimal transport and sliced optimal transport.

2.1 Notations

In what follows, $\mathcal{M}_+(\mathbb{R}^d)$ denotes the set of all positive Radon measures of finite mass on \mathbb{R}^d . For any $\alpha \in \mathcal{M}_+(\mathbb{R}^d)$, $\text{supp}(\alpha)$ is the support of α , and $m(\alpha) = \int_{\mathbb{R}^d} d\alpha(x) < +\infty$ is the mass of α . For $\alpha \in \mathcal{M}_+(\mathbb{R}^d)$ and a map $T : \mathbb{R}^d \rightarrow \mathbb{R}^p$, $T_{\#}\alpha$ is the pushforward measure of α by T , defined for all $A \subset \mathbb{R}^p$ as $T_{\#}\mu(A) = \mu(T^{-1}(A))$. Let δ_z be the Dirac measure at z and for $n \geq 1$, define the empirical measure $\hat{\alpha}_n = \sum_{i=1}^n w_i \delta_{Z_i}$, where $(Z_i)_{i=1}^n$ are n independent and identically distributed (i.i.d.) samples from $\alpha \in \mathcal{M}_+(\mathbb{R}^d)$, and $w_i > 0$. For any convex function $\varphi : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$, we denote by φ^* its Legendre transform, *i.e.*, for $x \in \mathbb{R}$, $\varphi^*(x) = \sup_{y \geq 0} xy - \varphi(y)$. We will also use the notation $\varphi^\circ(x) = -\varphi^*(-x)$. For $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^d)$, $\alpha \otimes \beta$ is the product measure, and for $f, g : \mathbb{R}^d \rightarrow \mathbb{R}$, denote by $f \oplus g : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ the mapping defined as $(f \oplus g)(x, y) = f(x) + g(y)$ for all $x, y \in \mathbb{R}^d$. $\mathbb{S}^{d-1} \triangleq \{\theta \in \mathbb{R}^d : \|\theta\| = 1\}$ is the unit sphere, and for $\theta \in \mathbb{S}^{d-1}$, $\theta^* : \mathbb{R}^d \rightarrow \mathbb{R}$ is the mapping defined as $\theta^*(x) = \langle \theta, x \rangle$ for all $x \in \mathbb{R}^d$.

2.2 Unbalanced Optimal Transport

We recall the static formulation of unbalanced OT proposed by Liero et al. (2018), which uses φ -divergences as penalty terms.

Definition 2.1 (φ -divergences). *Let $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$. Let $\varphi : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ be an entropy function, *i.e.*, φ is convex, lower semicontinuous, $\text{dom}(\varphi) \triangleq \{x \in \mathbb{R}, \varphi(x) < +\infty\} \subset [0, +\infty)$ and $\varphi(1) = 0$. Denote $\varphi'_\infty \triangleq \lim_{x \rightarrow +\infty} \frac{\varphi(x)}{x}$. The φ -divergence between α and β is*

$$D_\varphi(\alpha|\beta) \triangleq \int_{\mathbb{R}^d} \varphi\left(\frac{d\alpha}{d\beta}(x)\right) d\beta(x) + \varphi'_\infty \int_{\mathbb{R}^d} d\alpha^\perp(x), \quad (1)$$

where α^\perp is defined as $\alpha = \frac{d\alpha}{d\beta}\beta + \alpha^\perp$.

Definition 2.2 (Unbalanced OT (Liero et al., 2018)). *Let (φ_1, φ_2) be a pair of entropy functions and $C_d : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ a cost function. The unbalanced OT problem between $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$ reads*

$$\text{UOT}(\alpha, \beta) \triangleq \inf_{\pi \in \mathcal{M}_+(\mathbb{R}^d \times \mathbb{R}^d)} \int C_d(x, y) d\pi(x, y) + D_{\varphi_1}(\pi_1|\alpha) + D_{\varphi_2}(\pi_2|\beta), \quad (2)$$

where π_1 and π_2 denote the marginal distributions of π with respect to (*w.r.t.*) the first and second variable respectively.

When $\varphi_1 = \varphi_2$ and $\varphi_1(x) = 0$ for $x = 1$, $\varphi_1(x) = +\infty$ otherwise, (2) boils down to the Kantorovich formulation of OT (or *balanced OT*), denoted by $\text{OT}(\alpha, \beta)$. Indeed, in that case, $D_{\varphi_1}(\pi_1|\alpha) = D_{\varphi_2}(\pi_2|\beta) = 0$ if $\pi_1 = \alpha$ and $\pi_2 = \beta$, $D_{\varphi_1}(\pi_1|\alpha) = D_{\varphi_2}(\pi_2|\beta) = +\infty$ otherwise.

Under other suitable choices of entropy functions φ_1 and φ_2 , $\text{UOT}(\alpha, \beta)$ is more robust than $\text{OT}(\alpha, \beta)$, since it can discard outliers and compare α and β with different masses. We refer to (Séjourné et al., 2022a, Section 4.2) for a detailed discussion on the choice of entropies and its consequences on the transport plan computed by UOT. Two common choices are $\varphi_i(x) = \rho|x - 1|$ and $\varphi_i(x) = \rho(x \log(x) - x + 1)$, where $\rho > 0$ is a characteristic radius *w.r.t.* C_d . They respectively correspond to $D_{\varphi_i} = \rho\text{TV}$ (total variation distance (Chizat et al., 2018a)) and $D_{\varphi_i} = \rho\text{KL}$ (Kullback-Leibler divergence), and operate differently: KL smooths out geometric outliers, while TV either keeps or removes samples (Séjourné et al., 2022a). The GHK distance corresponds to (2) with $C_d(x, y) = \|x - y\|^2$ and $D_{\varphi_i} = \rho_i\text{KL}$ (Liero et al., 2018).

One can obtain an equivalent formulation of UOT by deriving the dual of (2) and proving strong duality. We recall this result below.

Proposition 2.3 (Corollary 4.12 in (Liero et al., 2018)). *The UOT problem (2) can equivalently be written as $\text{UOT}(\alpha, \beta) = \sup_{f \oplus g \leq C_d} \mathcal{D}(f, g; \alpha, \beta)$, with*

$$\mathcal{D}(f, g; \alpha, \beta) \triangleq \int \varphi_1^\circ(f(x)) d\alpha(x) + \int \varphi_2^\circ(g(y)) d\beta(y), \quad (3)$$

where for $i \in \{1, 2\}$, $\varphi_i^\circ(x) \triangleq -\varphi_i^*(-x)$ with $\varphi_i^*(x) \triangleq \sup_{y \geq 0} xy - \varphi_i(y)$ the Legendre transform of φ_i , and $f \oplus g \leq C_d$ means that for $(x, y) \sim \alpha \otimes \beta$, $f(x) + g(y) \leq C_d(x, y)$.

When clear from the context, we will omit the dependence on (α, β) and write $\mathcal{D}(f, g)$ instead of $\mathcal{D}(f, g; \alpha, \beta)$. The Legendre transform of φ_i is well known for typical choices of φ_i -divergences. For example, if $D_{\varphi_i} = \rho_i \text{KL}$, then $\varphi_i^*(x) = \rho_i(e^{x/\rho_i} - 1)$.

Based on Proposition 2.3, one can compute $\text{UOT}(\alpha, \beta)$ by optimizing a pair of continuous functions (f, g) . However, $\text{UOT}(\alpha, \beta)$ is known to be computationally intensive (Pham et al., 2020), which motivates the development of methods able to scale to the large dimensions and sample sizes encountered in ML applications.

2.3 Sliced Optimal Transport

Among the many workarounds that have been proposed to overcome the OT computational bottleneck (Peyré et al., 2019), Sliced OT (Rabin et al., 2012) has attracted a lot of attention due to its computational benefits and theoretical guarantees.

Definition 2.4 (Sliced OT). *Let $\mathbb{S}^{d-1} \triangleq \{\theta \in \mathbb{R}^d : \|\theta\| = 1\}$ be the unit sphere in \mathbb{R}^d . For $\theta \in \mathbb{S}^{d-1}$, denote by $\theta^* : \mathbb{R}^d \rightarrow \mathbb{R}$ the linear map such that for $x \in \mathbb{R}^d$, $\theta^*(x) \triangleq \langle \theta, x \rangle$. Let σ be the uniform probability over \mathbb{S}^{d-1} . Consider $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$. The Sliced OT problem is defined as*

$$\text{SOT}(\alpha, \beta) \triangleq \int_{\mathbb{S}^{d-1}} \text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta) d\sigma(\theta), \quad (4)$$

where for any measurable function f and $\xi \in \mathcal{M}_+(\mathbb{R}^d)$, $f_\# \xi$ is the push-forward measure of ξ by f , i.e., for any measurable set $A \subset \mathbb{R}$, $f_\# \xi(A) \triangleq \xi(f^{-1}(A))$, $f^{-1}(A) \triangleq \{x \in \mathbb{R}^d : f(x) \in A\}$.

Since $(\theta_\#^* \alpha, \theta_\#^* \beta)$ are two measures supported on \mathbb{R} , $\text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ is defined in terms of a cost function $C_1 : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, and can be efficiently computed. Therefore, $\text{SOT}(\alpha, \beta)$ can provide significant computational advantages over $\text{OT}(\alpha, \beta)$ in large-scale settings. In practice, if $\alpha = \sum_{i=1}^n \alpha_i \delta_{x_i}$ and $\beta = \sum_{i=1}^n \beta_i \delta_{y_i}$ are discrete measures, the standard procedure for approximating $\text{SOT}(\alpha, \beta)$ consists in sampling m i.i.d. samples $\{\theta_j\}_{j=1}^m$ from σ , then computing $\text{OT}((\theta_j^*)_\# \alpha, (\theta_j^*)_\# \beta)$ for $j = 1, \dots, m$. This second step involves sorting the n support points of α and β (Peyré et al., 2019, Section 2.6), thus involves $\mathcal{O}(n \log n)$ operations per θ_j .

$\text{SOT}(\alpha, \beta)$ relies on the Kantorovich formulation of OT, thus $\text{SOT}(\alpha, \beta) < +\infty$ only when $m(\alpha) = m(\beta)$, and may not provide meaningful comparisons in presence of outliers. To overcome such limitations, prior works have proposed slicing a particular instance of UOT that is partial OT (Bonneel & Coeurjolly, 2019; Bai et al., 2023), for which D_φ is the total variation distance. More precisely, noting POT the UOT problem with $D_\varphi = \rho \text{TV}$, they consider for $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^d)$ the problem

$$\text{SPOT}(\alpha, \beta) \triangleq \int_{\mathbb{S}^{d-1}} \text{POT}(\theta_\#^* \alpha, \theta_\#^* \beta) d\sigma(\theta). \quad (5)$$

For the 1D partial OT problem, Bonneel & Coeurjolly (2019) solve a one dimensional injective partial assignment in quasilinear complexity, but which does not allow for mass destruction in the source measure, while Bai et al. (2023) proposed an efficient procedure with a quadratic worst case complexity. However, their algorithms only apply to measures whose samples have constant mass (e.g., $\alpha_i = \beta_j = 1$). In the next section, we generalize their line of work and propose a new way of combining sliced OT and unbalanced OT.

3 Sliced Unbalanced OT and Unbalanced Sliced OT

We present two new scalable and robust OT problems, by combining the unbalanced and slicing strategies in two different ways. We conduct a theoretical analysis of both strategies and provide a comparison of the two. For ease of exposition, all proofs of the results in this section are provided in Appendix A.

First, we propose to *slice the unbalanced OT problem*: we average the UOT problem over different projections of the compared measures, similar to the approach of sliced partial OT (Bonnel & Coeurjolly, 2019; Bai et al., 2023). We refer to this problem as Sliced Unbalanced OT (SUOT) and introduce it in Section 3.1. Next, we explore the reverse strategy, *i.e.*, we *unbalance the sliced OT problem*: the weights of SUOT are penalized to introduce imbalance, analogous to how UOT relates to OT. We call this method *Unbalanced Sliced OT* (USOT) and present it in Section 3.2.

3.1 Sliced Unbalanced Optimal Transport

Our first strategy consists in slicing the unbalanced OT problem and leads to the following definition.

Definition 3.1 (Sliced Unbalanced OT). *Let (φ_1, φ_2) be a pair of entropy functions and $C_1 : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ a cost function. The sliced unbalanced OT problem between $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$ reads*

$$\text{SUOT}(\alpha, \beta) \triangleq \int_{\mathbb{S}^{d-1}} \text{UOT}(\theta_{\sharp}^* \alpha, \theta_{\sharp}^* \beta) d\sigma(\theta), \quad (6)$$

where for $\theta \sim \sigma$, $\text{UOT}(\theta_{\sharp}^* \alpha, \theta_{\sharp}^* \beta) = \inf_{\pi_{\theta} \in \mathcal{M}_+(\mathbb{R} \times \mathbb{R})} \int_{\mathbb{R} \times \mathbb{R}} C_1(x, y) d\pi_{\theta}(x, y) + D_{\varphi_1}((\pi_{\theta})_1 | \theta_{\sharp}^* \alpha) + D_{\varphi_2}((\pi_{\theta})_2 | \theta_{\sharp}^* \beta)$ with $(\pi_{\theta})_1, (\pi_{\theta})_2$ the marginal distributions of π_{θ} .

By definition, SUOT is a specific instance of the class of sliced probability divergences (Nadjahi et al., 2020a), where the *base divergence* is chosen as UOT. SUOT can also be interpreted as a general expression of the sliced partial OT problem (Bonnel & Coeurjolly, 2019; Bai et al., 2023): while the latter imposes $D_{\varphi_i} = \rho_i \text{TV}$, SUOT allows for the use of arbitrary φ -divergences.

In the following, we establish a set of theoretical properties for SUOT with different choices of φ -divergences and cost functions C_1 . First, we identify sufficient conditions for which the solution of (6) exists.

Proposition 3.2 (SUOT: Existence of solutions). *Assume that C_1 is lower-semicontinuous and that either (i) $\varphi'_{1,\infty} = \varphi'_{2,\infty} = +\infty$, or (ii) C_1 has compact sublevels on $\mathbb{R} \times \mathbb{R}$ and $\varphi'_{1,\infty} + \varphi'_{2,\infty} + \inf C_1 > 0$. Then, the solution of $\text{SUOT}(\alpha, \beta)$ exists, in the sense that for any $\theta \sim \sigma$, there exists $\pi_{\theta}^* \in \mathcal{M}_+(\mathbb{R} \times \mathbb{R})$ attaining the infimum in $\text{UOT}(\theta_{\sharp}^* \alpha, \theta_{\sharp}^* \beta)$.*

The assumptions of Proposition 3.2 are met for some settings of interest, including $D_{\varphi_1} = D_{\varphi_2} = \text{KL}$ (since $\varphi'_{\infty} = +\infty$), or $D_{\varphi_1} = D_{\varphi_2} = \text{TV}$ and $C_1(x, y) = |x - y|^p$ ($p \geq 1$) (since $\varphi'_{\infty} = 1$): see (Séjourné et al., 2022a, Section 2.1) for more details.

Next, we show some topological properties of SUOT. In the next proposition, we prove that SUOT preserves the metric properties of UOT, which is consistent with (Nadjahi et al., 2020a, Proposition 1). In Section 3.3, we study the metrization of the weak*-topology with SUOT.

Proposition 3.3 (SUOT: Metric properties). *Suppose UOT is non-negative, symmetric and/or definite on $\mathcal{M}_+(\mathbb{R}) \times \mathcal{M}_+(\mathbb{R})$. Then, SUOT is respectively non-negative, symmetric and/or definite on $\mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$. If there exists $p \in [1, +\infty)$ s.t. for any $\alpha, \beta, \gamma \in \mathcal{M}_+(\mathbb{R})$, $\text{UOT}^{1/p}(\alpha, \beta) \leq \text{UOT}^{1/p}(\alpha, \gamma) + \text{UOT}^{1/p}(\gamma, \beta)$, then $\text{SUOT}^{1/p}(\alpha, \beta) \leq \text{SUOT}^{1/p}(\alpha, \gamma) + \text{SUOT}^{1/p}(\gamma, \beta)$.*

By Proposition 3.3, establishing the metric axioms of UOT between *univariate* measures (as detailed in (Séjourné et al., 2022a, Section 3.3.1)) is sufficient to prove the metric properties of SUOT between *multivariate* measures. For example, since GHK is a metric for the order $p = 2$ (Liero et al., 2018), so is the induced SUOT.

We move on to the statistical aspects and study the sample complexity of SUOT. For $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$, we establish the speed of convergence of $\text{SUOT}(\hat{\alpha}_n, \hat{\beta}_n)$ toward $\text{SUOT}(\alpha, \beta)$, where $\hat{\alpha}_n, \hat{\beta}_n$ denote

the empirical measures supported on n independent samples from α, β respectively (as defined in Section 2.1). We prove below that SUOT extends the sample complexity of UOT from one-dimensional settings to multi-dimensional ones.

Theorem 3.4 (SUOT: Sample complexity). (i) *Assume for $(\mu, \nu) \in \mathbb{M} \times \mathbb{M}$ with $\mathbb{M} \subset \mathcal{M}_+(\mathbb{R})$, $\mathbb{E}|\text{UOT}(\mu, \nu) - \text{UOT}(\hat{\mu}_n, \hat{\nu}_n)| \leq \kappa(n)$. Then, for $(\alpha, \beta) \in \tilde{\mathbb{M}} \times \tilde{\mathbb{M}}$ with $\tilde{\mathbb{M}} \triangleq \{\eta \in \mathcal{M}_+(\mathbb{R}^d) : \forall \theta \in \mathbb{S}^{d-1}, \theta_{\#}^* \eta \in \mathbb{M}\}$, $\mathbb{E}|\text{SUOT}(\alpha, \beta) - \text{SUOT}(\hat{\alpha}_n, \hat{\beta}_n)| \leq \kappa(n)$.*

(ii) *Assume for $\mu \in \mathbb{M}$ with $\mathbb{M} \subset \mathcal{M}_+(\mathbb{R})$, $\mathbb{E}|\text{UOT}(\mu, \hat{\mu}_n)| \leq \xi(n)$. Then, for $\alpha \in \tilde{\mathbb{M}}$ with $\tilde{\mathbb{M}} \triangleq \{\eta \in \mathcal{M}_+(\mathbb{R}^d) : \forall \theta \in \mathbb{S}^{d-1}, \theta_{\#}^* \eta \in \mathbb{M}\}$, $\mathbb{E}|\text{SUOT}(\alpha, \hat{\alpha}_n)| \leq \xi(n)$.*

Note that the expectations in Theorem 3.4 are taken with respect to the samples of the empirical measures, which are random. Theorem 3.4 shows that SUOT enjoys a *dimension-free* sample complexity, even when comparing multivariate measures. This advantage is recurrent of sliced divergences (Nadjahi et al., 2020b) and further motivates their use on high-dimensional settings. The sample complexity rates $\kappa(n)$ or $\xi(n)$ can be deduced from the literature on UOT for univariate measures. For instance, in the GHK setting, the rate is given by $\kappa(n) \propto n^{-1/2}$ for measures with compact, convex support and continuously differentiable densities (Vacher & Vialard, 2023, Corollary 3.4), and a suitable class \mathbb{M} can be defined.

Finally, we derive the dual formulation of SUOT and prove that strong duality holds. This result has important practical implications, as we will leverage it in Section 4 to develop a methodology for computing SUOT. Note that the computation of SUOT involves integration with respect to σ , which generally cannot be done in closed form, as is the case for most sliced divergences. Since our goal is to develop a practical and implementable method, we will consider the Monte Carlo approximation commonly used by practitioners to compute sliced divergences (Nadjahi et al., 2020a): we approximate $\text{SUOT}(\alpha, \beta)$ as $\int_{\mathbb{S}^{d-1}} \text{UOT}(\theta_{\#}^* \alpha, \theta_{\#}^* \beta) d\hat{\sigma}_K(\theta)$, where $\hat{\sigma}_K$ is a discrete distribution supported on K i.i.d. samples drawn from σ .

Theorem 3.5 (SUOT: Strong duality). *For $i \in \{1, 2\}$, let φ_i be an entropy function such that $\text{dom}(\varphi_i^*) \cap \mathbb{R}_-$ is non-empty, and either $0 \in \text{dom}(\varphi_i)$ or $m(\alpha), m(\beta) \in \text{dom}(\varphi_i)$. Let $\mathcal{E} \triangleq \{(f_\theta, g_\theta)_{\theta \in \text{supp}(\hat{\sigma}_K)} : \forall \theta \in \text{supp}(\hat{\sigma}_K), f_\theta \oplus g_\theta \leq C_1\}$. Then,*

$$\text{SUOT}(\alpha, \beta) = \sup_{(f_\theta, g_\theta) \in \mathcal{E}} \int_{\mathbb{S}^{d-1}} \mathcal{D}(f_\theta, g_\theta; \theta_{\#}^* \alpha, \theta_{\#}^* \beta) d\hat{\sigma}_K(\theta). \quad (7)$$

3.2 Unbalanced Sliced Optimal Transport

As a second strategy to make unbalanced OT scalable, we propose to unbalance sliced OT. To this end, we start with the following formulation of UOT (Liero et al., 2018),

$$\text{UOT}(\alpha, \beta) = \inf_{(\pi_1, \pi_2) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)} \text{OT}(\pi_1, \pi_2) + \text{D}_{\varphi_1}(\pi_1 | \alpha) + \text{D}_{\varphi_2}(\pi_2 | \beta), \quad (8)$$

and we replace UOT by its sliced counterpart, SOT. This yields the following definition:

Definition 3.6 (Unbalanced Sliced OT). *Let (φ_1, φ_2) be a pair of entropy functions. The unbalanced sliced OT problem between $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$ reads*

$$\text{USOT}(\alpha, \beta) \triangleq \inf_{(\pi_1, \pi_2) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)} \text{SOT}(\pi_1, \pi_2) + \text{D}_{\varphi_1}(\pi_1 | \alpha) + \text{D}_{\varphi_2}(\pi_2 | \beta). \quad (9)$$

This approach is entirely novel since, to the best of our knowledge, it has never been studied in prior work, even for specific choices of entropy. To gain a better grasp of this new object, USOT, we examine how the theoretical properties discussed in the previous section apply here.

We first prove that the solution of (9) exists under the same conditions as those for SUOT outlined in Proposition 3.2.

Proposition 3.7 (USOT: Existence of solutions). *Assume that C_1 is lower-semicontinuous and that either (i) $\varphi'_{1,\infty} = \varphi'_{2,\infty} = +\infty$, or (ii) C_1 has compact sublevels on $\mathbb{R} \times \mathbb{R}$ and $\varphi'_{1,\infty} + \varphi'_{2,\infty} + \inf C_1 > 0$. Then, the solution of $\text{USOT}(\alpha, \beta)$ exists: there exists $(\pi_1^*, \pi_2^*) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$ attaining the infimum in (9).*

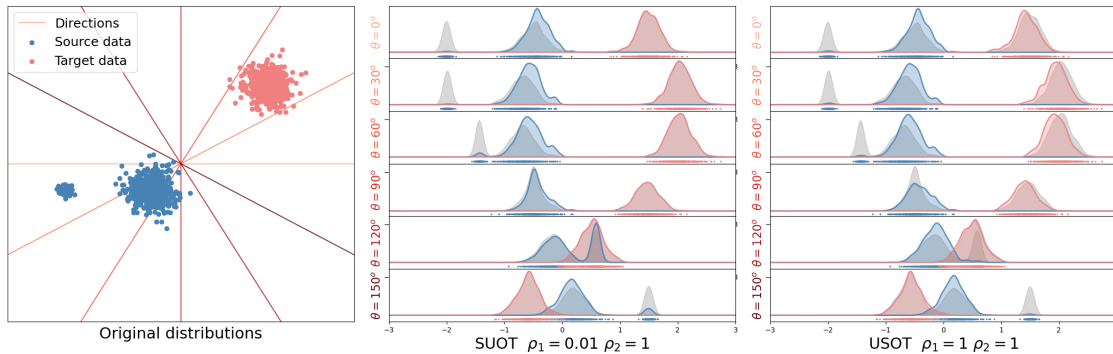


Figure 1: **Toy illustration** on the behaviors of SUOT and USOT. (left) Original 2D samples and slices used for illustration. KDE density estimations of the projected samples: grey, original distributions, colored, distributions reweighed by SUOT (center), and reweighed by USOT (right).

We then review the conditions under which USOT is a (pseudo-)metric, and we prove that strong duality holds. Similar to SUOT, the dual formulation that we derive will enable the design of an algorithm for effectively computing USOT in practice.

Proposition 3.8 (USOT: Metric properties). *For any $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$, $\text{USOT}(\alpha, \beta) \geq 0$. If $\varphi_1 = \varphi_2$, USOT is symmetric. If D_{φ_1} and D_{φ_2} are definite, then USOT is definite. If $C_1(x, y) = |x - y|$ and $D_{\varphi_1} = D_{\varphi_2} = \rho \text{TV}$, then, USOT satisfies the triangle inequality.*

Theorem 3.9 (USOT: Strong duality). *For $i \in \{1, 2\}$, let φ_i be an entropy function such that $\text{dom}(\varphi_i^*) \cap \mathbb{R}_-$ is non-empty, and either $0 \in \text{dom}(\varphi_i)$ or $m(\alpha), m(\beta) \in \text{dom}(\varphi_i)$. Let $\mathcal{E} \triangleq \{(f_\theta, g_\theta)_{\theta \in \text{supp}(\hat{\sigma}_K)} : \forall \theta \in \text{supp}(\hat{\sigma}_K), f_\theta \oplus g_\theta \leq C_1\}$. Then,*

$$\text{USOT}(\alpha, \beta) = \sup_{(f_\theta, g_\theta) \in \mathcal{E}} \mathcal{D} \left(\int_{\mathbb{S}^{d-1}} f_\theta \circ \theta^* d\hat{\sigma}_K(\theta), \int_{\mathbb{S}^{d-1}} g_\theta \circ \theta^* d\hat{\sigma}_K(\theta); \alpha, \beta \right). \quad (10)$$

Since USOT does not belong to the class of sliced divergences, establishing its sample complexity is more challenging compared to SUOT. Based on the literature, one standard technique involves deriving covering number bounds on the space of the dual potentials of USOT. This theoretical question is highly non-trivial given the complex structure of \mathcal{E} , and as such is out of the scope of this paper. Nevertheless, we investigate the sample complexity on empirical settings: our experimental results presented in Appendix C.5 suggest that USOT might also enjoy a dimension-free rate.

3.3 Comparative Analysis of Sliced Unbalanced and Unbalanced Sliced Optimal Transport

In addition to the theoretical analysis previously conducted for SUOT and USOT independently, this section provides further insights to better grasp the differences between these two strategies.

First, by comparing Definition 3.1 with Definition 3.6, SUOT and USOT clearly differ at the conceptual level. Specifically, $\text{SUOT}(\alpha, \beta)$ penalizes the marginals of π_θ for $\theta \sim \sigma$, where π_θ is the coupling that transports mass from $\theta_\#^* \alpha$ to $\theta_\#^* \beta$. In contrast, $\text{USOT}(\alpha, \beta)$ directly regularizes the marginals of the coupling between α and β . To illustrate this difference, we consider $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^2) \times \mathcal{M}_+(\mathbb{R}^2)$ with α contaminated by outliers, then compute $\text{SUOT}(\alpha, \beta)$ and $\text{USOT}(\alpha, \beta)$. We plot (α, β) and the sampled projections $(\theta_k)_k$ (Figure 1, left), the marginals of $(\pi_{\theta_k})_k$ obtained with $\text{SUOT}(\alpha, \beta)$ (Figure 1, center), and the marginals of $((\theta_k^*)_\# \pi)_k$ with $\text{USOT}(\alpha, \beta)$ (Figure 1, right). We observe that the source outliers in α have been successfully removed by $\text{USOT}(\alpha, \beta)$ for all θ_k , while they may still appear with $\text{SUOT}(\alpha, \beta)$ (e.g., Figure 1, center: note the bimodal marginal in blue for $\theta = 120^\circ$). This difference is due to the marginal penalization terms in $\text{USOT}(\alpha, \beta)$, which operate directly w.r.t. (α, β) rather than their projections $(\theta_\#^* \alpha, \theta_\#^* \beta)$, unlike $\text{SUOT}(\alpha, \beta)$.

A question of particular interest regarding probability divergences is how they relate to each other, specifically whether they yield equivalent topologies. We explore this question for SUOT and USOT. To do so, we

consider a notion of equivalence that is frequently studied in the literature, as in (Bayraktar & Guo, 2021, Theorem 2.3.(i)).

We start by proving a first set of inequalities that relate SUOT, USOT and UOT: our next theorem shows that USOT is always greater than SUOT and that, for appropriate choices of cost functions, USOT is upper-bounded by UOT.

Theorem 3.10. *For any $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$,*

$$\text{SUOT}(\alpha, \beta) \leq \text{USOT}(\alpha, \beta). \quad (11)$$

Moreover, suppose that, $\forall(x, y) \in \mathbb{R}^d \times \mathbb{R}^d, \forall \theta \in \mathbb{S}^{d-1}, C_1(\theta^(x), \theta^*(y)) \leq C_d(x, y)$. Then, for any $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$,*

$$\text{SUOT}(\alpha, \beta) \leq \text{USOT}(\alpha, \beta) \leq \text{UOT}(\alpha, \beta). \quad (12)$$

In particular, Theorem 3.10 holds for the following common choice of costs: $\forall(s, t) \in \mathbb{R}^2, C_1(s, t) = |s - t|^p$ and $\forall(x, y) \in \mathbb{R}^d \times \mathbb{R}^d, C_d(x, y) = \|x - y\|^p$, with $p \in [1, +\infty)$.

Next, we prove that $\text{UOT}(\alpha, \beta)$ can be upper-bounded by a functional of $\text{SUOT}(\alpha, \beta)$ when (α, β) have compact supports, by adapting the reasoning from Bonnotte (2013, Lemma 5.1.4) to our setting and considering the duals of UOT (Proposition 2.3) and SUOT (Theorem 3.5) instead of the dual of OT and SOT. Most arguments in (Bonnotte, 2013) adapt well to our setting, but establishing a Lipschitz condition on the integrand of the dual required a more technical approach. To this end, we prove Lemma A.13, which results in a different constant value, denoted as $c(m(\alpha), m(\beta), \rho, R)$.

Theorem 3.11. *Let $X \subset \mathbb{R}^d$ be a compact set with radius R . Define the cost functions as $C_1(s, t) = |s - t|^p, (s, t) \in \mathbb{R}^2$, and $C_d(x, y) = \|x - y\|^p, (x, y) \in \mathbb{R}^d \times \mathbb{R}^d$, with $p \in [1, +\infty)$. Assume either, (i) $D_{\varphi_1} = D_{\varphi_2} = \rho\text{KL}$; or (ii) $p = 1$ and $D_{\varphi_1} = D_{\varphi_2} = \rho\text{TV}$. Then, for any $(\alpha, \beta) \in \mathcal{M}_+(X) \times \mathcal{M}_+(X)$,*

$$\text{UOT}(\alpha, \beta) \leq c(m(\alpha), m(\beta), \rho, R) \text{SUOT}(\alpha, \beta)^{\frac{1}{a+1}},$$

where $c(m(\alpha), m(\beta), \rho, R)$ is a constant depending on $m(\alpha), m(\beta), \rho, R$, which is non-decreasing in $m(\alpha)$ and $m(\beta)$.

We show the equivalence of SUOT, USOT and UOT by combining Theorem 3.11 and Theorem 3.10, assuming that the constant $c(m(\alpha), m(\beta), \rho, R)$ does not depend on $m(\alpha), m(\beta)$. This occurs, for example, when the masses of α and β are uniformly bounded; that is, there exists $M \in \mathbb{R}_+$ such that $m(\alpha) \leq M$ and $m(\beta) \leq M$.

The equivalence of SUOT, USOT and UOT is a key result for proving that SUOT and USOT metrize weak* convergence, provided that UOT does (as in the GHK setting (Liero et al., 2018, Theorem 7.25)). Recall that a sequence of positive measures $(\alpha_n)_{n \in \mathbb{N}^*}$ converges weakly to $\alpha \in \mathcal{M}_+(\mathbb{R}^d)$ (denoted by $\alpha_n \rightharpoonup \alpha$) if, for any continuous and bounded $f : \mathbb{R}^d \rightarrow \mathbb{R}$, $\lim_{n \rightarrow +\infty} \int f d\alpha_n = \int f d\alpha$.

Theorem 3.12 (Metrizability of the weak* topology by SUOT, USOT). *Assume the conditions in Theorem 3.11 are met. Let $(\alpha_n)_{n \in \mathbb{N}^*}$ be a sequence of measures in $\mathcal{M}_+(X)$ and $\alpha \in \mathcal{M}_+(X)$, where $X \subset \mathbb{R}^d$ is a compact set with radius R . Then, SUOT and USOT metrize the weak* convergence, i.e., $\alpha_n \rightharpoonup \alpha \Leftrightarrow \lim_{n \rightarrow +\infty} \text{SUOT}(\alpha_n, \alpha) = 0$, and $\alpha_n \rightharpoonup \alpha \Leftrightarrow \lim_{n \rightarrow +\infty} \text{USOT}(\alpha_n, \alpha) = 0$.*

The metrizability of weak* convergence was not studied in related work, including in existing instances of our framework, such as partial OT (Bonnel & Coeurjolly, 2019; Bai et al., 2023). In addition to complementing prior work, our result paves the way for other research directions. For instance, it can be used to justify the well-posedness of approximating an unbalanced Wasserstein gradient flow (Ambrosio et al., 2005) using SUOT, as done for SOT in (Candau-Tilh, 2020; Bonet et al., 2022). Unbalanced Wasserstein gradient flows have been a key tool in deep learning theory, e.g., to prove global convergence of one-hidden layer neural networks (Chizat & Bach, 2018; Rotskoff et al., 2019).

4 Computing SUOT and USOT with Frank-Wolfe algorithms

In this section, we propose two algorithms to compute SUOT and USOT in practice. The resulting procedures are given in Algorithms 2 and 3 respectively, and require smooth penalty terms $(D_{\varphi_1}, D_{\varphi_2})$. This condition

is satisfied in the GHK setting ($D_{\varphi_i} = \rho_i \text{KL}$), but not for sliced partial OT ($D_{\varphi_i} = \rho_i \text{TV}$, Bai et al. (2023)). Our strategy is inspired by Séjourné et al. (2022b), where they proposed to solve the unbalanced OT problem between univariate measures using the Frank-Wolfe algorithm (as recalled in Appendix B.2). More precisely, we apply FW to optimize translation-invariant forms of the dual problems derived in Theorems 3.5 and 3.9.

4.1 Background: Frank-Wolfe Algorithm and Application to One-Dimensional Unbalanced OT

FW is a popular iterative first-order optimization algorithm for solving $\max_{x \in \mathcal{E}} \mathcal{H}(x)$, where \mathcal{E} is a compact convex set and $\mathcal{H} : \mathcal{E} \rightarrow \mathbb{R}$ a concave, differentiable function. The procedure consists in maximizing a linear approximation of \mathcal{H} at each iteration: given the current iterate x_t , FW solves the *linear oracle* $r_{t+1} \in \arg \max_{r \in \mathcal{E}} \langle \nabla \mathcal{H}(x_t), r \rangle$, then performs $x_{t+1} = (1 - \gamma_{t+1})x_t + \gamma_{t+1}r_{t+1}$ with stepsize γ_{t+1} typically chosen as $\gamma_{t+1} = \frac{2}{2+t+1}$. We refer to this step as **FWStep** and report the pseudo-code in Appendix B.2.

Séjourné et al. (2022b) apply FW to solve a translation-invariant formulation of the dual of $\text{UOT}(\alpha, \beta)$ for $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}) \times \mathcal{M}_+(\mathbb{R})$, and show that the linear oracle in **FWStep** is the dual of $\text{OT}(\alpha_t, \beta_t)$ where (α_t, β_t) are normalized versions of (α, β) , i.e., $m(\alpha_t) = m(\beta_t) = 1$. Therefore, computing UOT amounts to solve a sequence of OT problems, which can efficiently be done since (α_t, β_t) are univariate probability measures. The expression of (α_t, β_t) depend on the input measures (α, β) , the current iterates (f_t, g_t) and the penalty coefficients (ρ_1, ρ_2) .

4.2 Frank-Wolfe Solvers for Sliced Unbalanced and Unbalanced Sliced OT

Translation-invariant duals. We compute $\text{SUOT}(\alpha, \beta)$ and $\text{USOT}(\alpha, \beta)$ for any $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$ by solving translation-invariant formulations of their duals with FW. By Theorems 3.5 and 3.9, and adapting the reasoning of Séjourné et al. (2022b), we prove that

$$\text{SUOT}(\alpha, \beta) = \sup_{(f_\theta, g_\theta) \in \mathcal{E}} \int_{\mathbb{S}^{d-1}} \mathcal{H}(f_\theta, g_\theta; \theta_{\#}^* \alpha, \theta_{\#}^* \beta) d\hat{\sigma}_K(\theta), \quad (13)$$

$$\text{USOT}(\alpha, \beta) = \sup_{(f_\theta, g_\theta) \in \mathcal{E}} \mathcal{H} \left(\int_{\mathbb{S}^{d-1}} f_\theta \circ \theta^* d\hat{\sigma}_K(\theta), \int_{\mathbb{S}^{d-1}} g_\theta \circ \theta^* d\hat{\sigma}_K(\theta); \alpha, \beta \right) \quad (14)$$

where $\mathcal{H}(f, g; \alpha, \beta) \triangleq \sup_{\lambda \in \mathbb{R}} \mathcal{D}(f + \lambda, g - \lambda; \alpha, \beta)$. These alternative duals are translation-invariant since, for any $\lambda \in \mathbb{R}$, $\mathcal{H}(f + \lambda, g - \lambda; \alpha, \beta) = \mathcal{H}(f, g; \alpha, \beta)$. If $(\varphi_1^\circ, \varphi_2^\circ)$ are smooth and strictly concave, then the maximizer in \mathcal{H} , denoted by $\lambda^*(f, g)$, exists and is unique. In particular, when $D_{\varphi_1} = \rho_1 \text{KL}$ and $D_{\varphi_2} = \rho_2 \text{KL}$, $\lambda^*(f, g)$ admits an analytical expression, which is given in the normalization routine (Algorithm 1). This is convenient as it avoids the need for approximate solvers to compute $\mathcal{H}(f, g; \alpha, \beta)$.

Frank-Wolfe iterations. We then apply FW to solve (13) and (14). We show that each iteration consists in solving a particular sliced OT problem between probability measures that depend on the input (α, β) and the iterates. To clarify this point, we present below the updates of **FWStep** tailored for each problem, starting with SUOT.

Proposition 4.1 (Frank-Wolfe iterations for SUOT). *Let $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$ and consider solving (13) with FW. Assume that $(\varphi_1^\circ, \varphi_2^\circ)$ are smooth and strictly concave. Given current iterates $(f_\theta^t, g_\theta^t)_{\theta \in \text{supp}(\hat{\sigma}_K)} \in \mathcal{E}$, the solutions of the linear oracle $(r_\theta^t, s_\theta^t)_{\theta \in \text{supp}(\hat{\sigma}_K)}$ are the dual potentials of $\int_{\mathbb{S}^{d-1}} \text{OT}(\theta_{\#}^* \alpha_\theta^t, \theta_{\#}^* \beta_\theta^t) d\hat{\sigma}_K(\theta)$, where $(\alpha_\theta^t, \beta_\theta^t)$ are measures given by $\alpha_\theta^t = \nabla \varphi_1^\circ(f_\theta^t + \lambda^*(f_\theta^t, g_\theta^t))\alpha$ and $\beta_\theta^t = \nabla \varphi_2^\circ(g_\theta^t - \lambda^*(f_\theta^t, g_\theta^t))\beta$.*

Proposition 4.1 shows that each FW iteration for solving the translation-invariant dual of $\text{SUOT}(\alpha, \beta)$ reduces to solving a balanced sliced OT problem: by (Séjourné et al., 2022b, Proposition 1), the measures $(\alpha_\theta^t, \beta_\theta^t)$ have the same mass, i.e., $m(\alpha_\theta^t) = m(\beta_\theta^t)$. When using KL-based penalty terms, the procedure for computing $(\alpha_\theta^t, \beta_\theta^t)$ is detailed in Algorithm 1, and reports the closed-form expression of $\lambda^*(f_\theta^t, g_\theta^t)$.

Algorithm 1 – Norm $(\alpha, \beta, f, g, \rho_1, \rho_2)$

Input: $\alpha, \beta, f, g, \rho_1, \rho_2$

Output: Normalized measures $(\bar{\alpha}, \bar{\beta})$

$$\lambda^* \leftarrow \frac{\rho_1 \rho_2}{\rho_1 + \rho_2} \log \left(\frac{\int e^{-f(x)/\rho_1} d\alpha(x)}{\int e^{-g(y)/\rho_2} d\beta(y)} \right)$$

$$\bar{\alpha} \leftarrow e^{-\frac{(f(x) + \lambda^*)}{\rho_1}} \alpha$$

$$\bar{\beta} \leftarrow e^{-\frac{(g(y) - \lambda^*)}{\rho_2}} \beta$$

Return $(\bar{\alpha}, \bar{\beta})$

Algorithm 2 – SUOT

Input: $\alpha, \beta, F, (\theta_k)_{k=1}^K, \rho_1, \rho_2$
Output: $\text{SUOT}(\alpha, \beta), (f_\theta, g_\theta)$

```

 $(f_\theta, g_\theta) \leftarrow (0, 0)$ 
for  $t = 0, 1, \dots, F - 1$  do
  for  $\theta \in (\theta_k)_{k=1}^K$  do
     $(\alpha_\theta, \beta_\theta) \leftarrow \text{Norm}(\theta_\#^* \alpha, \theta_\#^* \beta, f_\theta, g_\theta, \rho_1, \rho_2)$ 
     $(r_\theta, s_\theta) \leftarrow \text{SlicedDual}(\alpha_\theta, \beta_\theta)$ 
     $f_\theta \leftarrow (1 - \gamma_t) f_\theta + \gamma_t r_\theta$  (FWStep)
     $g_\theta \leftarrow (1 - \gamma_t) g_\theta + \gamma_t s_\theta$  (FWStep)
  end for
end for
Return  $\text{SUOT}(\alpha, \beta), (f_\theta, g_\theta)$  as in (7)

```

Algorithm 3 – USOT

Input: $\alpha, \beta, F, (\theta_k)_{k=1}^K, \rho_1, \rho_2$
Output: $\text{USOT}(\alpha, \beta), (f_{avg}, g_{avg})$

```

 $(f_\theta, g_\theta, f_{avg}, g_{avg}) \leftarrow (0, 0, 0, 0)$ 
for  $t = 0, 1, \dots, F - 1$  do
  for  $\theta \in (\theta_k)_{k=1}^K$  do
     $(\pi_1, \pi_2) \leftarrow \text{Norm}(\alpha, \beta, f_{avg}, g_{avg}, \rho_1, \rho_2)$ 
     $(r_\theta, s_\theta) \leftarrow \text{SlicedDual}(\theta_\#^* \pi_1, \theta_\#^* \pi_2)$ 
  end for
   $(r_{avg}, s_{avg}) \leftarrow \frac{1}{K} \sum_{k=1}^K r_{\theta_k}, \frac{1}{K} \sum_{k=1}^K s_{\theta_k}$ 
   $f_{avg} \leftarrow (1 - \gamma_t) f_{avg} + \gamma_t r_{avg}$  (FWStep)
   $g_{avg} \leftarrow (1 - \gamma_t) g_{avg} + \gamma_t s_{avg}$  (FWStep)
end for
Return  $\text{USOT}(\alpha, \beta), (f_{avg}, g_{avg})$  as in (10)

```

Each iteration requires computing the dual potentials of a sliced OT problem, which is non-trivial: previous implementations related to sliced OT only output the value of the loss, $\text{SOT}(\alpha, \beta)$, typically in the context of training generative models (Deshpande et al., 2019; Nguyen et al., 2020). We thus design two novel implementations in PyTorch (Paszke et al., 2019) to compute the dual potentials of sliced OT. The first one leverages that the gradient of $\text{OT}(\alpha, \beta)$ w.r.t. (α, β) are optimal (f, g) , which allows to backpropagate $\text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ w.r.t. (α, β) to obtain (r_θ, s_θ) . The second one computes them in parallel on GPUs using their closed form, which to the best of our knowledge, is a new sliced algorithm. We call $\text{SlicedDual}(\alpha, \beta)$ the step returning optimal (r_θ, s_θ) solving $\text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ for all $\theta \in \text{supp}(\hat{\sigma}_K)$, and refer to Appendix B.3 for the algorithms.

Building on Proposition 4.1 and the discussion above, we develop the FW methodology to compute $\text{SUOT}(\alpha, \beta)$ and detail it in Algorithm 2. Next, we derive the FW iterates for $\text{USOT}(\alpha, \beta)$.

Proposition 4.2 (Frank-Wolfe iterations for USOT). *Let $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$ and consider solving (14) with FW. Assume that $(\varphi_1^\circ, \varphi_2^\circ)$ are smooth and strictly concave. Given current iterates $(f_\theta^t, g_\theta^t)_{\theta \in \text{supp}(\hat{\sigma}_K)} \in \mathcal{E}$, the solutions of the linear oracle $(r_\theta^t, s_\theta^t)_{\theta \in \text{supp}(\hat{\sigma}_K)}$ are the dual potentials of $\text{SOT}(\bar{\alpha}^t, \bar{\beta}^t)$, where $(\bar{\alpha}_t, \bar{\beta}_t)$ are measures given by $\bar{\alpha}_t = \nabla \varphi^\circ(f_{avg} + \lambda^*(f_{avg}, g_{avg}))\alpha$ and $\bar{\beta}_t = \nabla \varphi^\circ(g_{avg} - \lambda^*(f_{avg}, g_{avg}))\beta$, with $f_{avg}(x) \triangleq \int_{\mathbb{S}^{d-1}} f_\theta^t(\theta^*(x)) d\hat{\sigma}_K(\theta)$, $g_{avg}(y) \triangleq \int_{\mathbb{S}^{d-1}} g_\theta^t(\theta^*(y)) d\hat{\sigma}_K(\theta)$.*

The resulting FW methodology, detailed in Algorithm 3, also leverages the Norm and SlicedDual routines. The key difference from $\text{SUOT}(\alpha, \beta)$ is in where the integral over $\theta \in \text{supp}(\hat{\sigma}_K)$ is performed, leading to a different balanced sliced OT problem to solve.

Marginals of UOT/USOT. The optimal primal marginals of UOT and USOT are geometric normalizations of inputs (α, β) with discarded outliers. Their computation involves the Norm routine detailed in Algorithm 1, using optimal dual potentials. This is how we compute marginals in Figure 1 and in the experiments of Section 5: see Appendix B.4 for more details.

4.3 Convergence Properties and Complexity

Convergence and stochastic Frank-Wolfe. Our theoretical setting verifies the assumptions of (Lacoste-Julien & Jaggi, 2015, Theorem 8), thus ensuring fast convergence of our methods. The number of FW iterations needed to converge remains low in our experiments. We give in Appendix B.5 empirical evidences that few iterations of FW ($F \leq 20$) suffice to reach numerical precision.

Formally, the preceding algorithms assume that the functional \mathcal{H} is given through integrals over the hypersphere, describing the set of all possible directions θ . However, in practice, SOT is computed by Monte-Carlo approximations, *i.e.*, drawing a fixed number K of directions $(\theta_k)_{k=1}^K$ and solving independently the different 1D OT problems. In the specific case of SUOT, this does not change much: K FW procedures are ran

Table 1: Accuracy on document classification

	BBCSport		Goodreads		
	Acc	t ($\cdot 10^{-3}$ s)	Acc (genre)	Acc (like)	t ($\cdot 10^{-3}$ s)
OT	94.55	3.12 \pm 1.61	55.22	71.00	440.30 \pm 250
UOT	96.73	243.39 \pm 9.24	-	-	-
SinkhUOT	95.45	46.22 \pm 2.17	53.55	67.81	2021.68 \pm 356
SOT	89.39 \pm 0.76	1.80 \pm 0.22	50.09 \pm 0.51	65.60 \pm 0.20	4.49 \pm 1.44
SUOT	90.12 \pm 0.15	13.9 \pm 1.21	50.15 \pm 0.04	66.72 \pm 0.38	14.32 \pm 0.95
USOT	93.52 \pm 0.04	14.37 \pm 1.29	52.67 \pm 0.62	67.78 \pm 0.39	14.45 \pm 0.88
SUOT (CV on ρ)	90.00 \pm 0.59	-	49.67 \pm 0.79	66.43 \pm 0.44	-
USOT (CV on ρ)	92.61 \pm 0.55	-	52.06 \pm 7.20	66.61 \pm 0.72	-

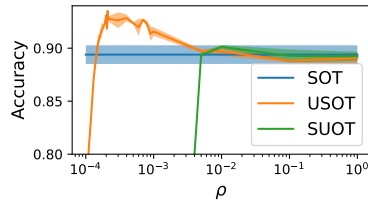


Figure 2: Ablation on BBCSport of ρ .

independently (eventually in parallel) over the fixed set of directions. The case of USOT relies on a global FW scheme, where f_{avg}, g_{avg} are computed w.r.t. a fixed distribution $\hat{\sigma}_K = (1/K) \sum_{k=1}^K \delta_{\theta_k}$. This empirical distribution of directions can be considered fixed throughout the FW iterations, or can be drawn independently for each iteration of the FW procedure. This actually corresponds to a *Stochastic FW* algorithm, which also converges as our setting verifies the assumptions of (Hazan & Luo, 2016, Theorem 3). We call this procedure *Stochastic USOT*, which corresponds to Algorithm 3 except that $(\theta_k)_{k=1}^K$ are sampled at each iteration. Since this procedure performs well in our experiments (e.g., Table 4) and $\mathbb{E}_{\theta_k \sim \sigma}[\hat{\sigma}_K] = \sigma$, this suggests the dual in Theorem 3.9 holds for σ .

Algorithmic complexity. FW algorithms and its variants have been widely studied theoretically. Computing SlicedDual has theoretically a complexity $\mathcal{O}(KN \log N)$, where N is the number of samples, and K the number of projections of $\hat{\sigma}_K$. However, we note that the sorting operation, which yields the super linear complexity, can be computed once for all FW iterations. Consequently, the overall complexity of SUOT and USOT is thus $\mathcal{O}(KN \log N + FN)$, where F is the number of FW iterations needed to reach convergence, with a $\mathcal{O}(N)$ complexity. Thus, our formulation enjoy a similar complexity than SOT, which is particularly appealing. However, *Stochastic USOT* is more costly, as each iteration requires sorting data projected along newly-sampled $(\theta_k)_{k=1}^K$. Its complexity is therefore $\mathcal{O}(KFN \log N)$. We finally note that due to the independent nature of the treatments of every projections, computing both Norm and SlicedDual operations can be done in parallel, leveraging GPU computations when available.

Extension to non-Euclidean settings. Interestingly, our algorithms offer great modularity, in the sense they can easily be used to compute unbalanced versions of existing variants of SOT. Indeed, while such variants differ in the one-dimensional representations of α and β they use, they all consist in solving 1D OT problems to compare α and β , which our FW strategy can solve. To illustrate this point, we combined our FW routine with hyperbolic SOT (Bonnet et al., 2023c) to compare measures supported on hyperbolic spaces: see Appendix C.3.

5 Experiments

This section presents a set of numerical experiments, which illustrate the effectiveness, robustness and computational efficiency of USOT¹. We first showcase the benefit of USOT over SUOT and SOT on a document classification task. Then, we consider experiments in very large scale settings such as color transfer on every pixels and the computation of barycenters of geophysical datasets.

5.1 Document classification

We first consider a document classification problem (Kusner et al., 2015). Documents are represented as distributions of words embedded with *word2vec* (Mikolov et al., 2013) in dimension $d = 300$. Let D_k be the k -th document and $x_1^k, \dots, x_{n_k}^k \in \mathbb{R}^d$ be the set of words in D_k . Then, $D_k = \sum_{i=1}^{n_k} w_i^k \delta_{x_i^k}$ where w_i^k is the frequency of x_i^k in D_k normalized s.t. $\sum_{i=1}^{n_k} w_i^k = 1$. Given a loss function L , the document classification task is solved by computing the matrix $(L(D_k, D_\ell))_{k,\ell}$, then using a k -nearest neighbor classifier. The aim

¹The code is available at https://github.com/clbonet/Slicing_Unbalanced_Optimal_Transport.

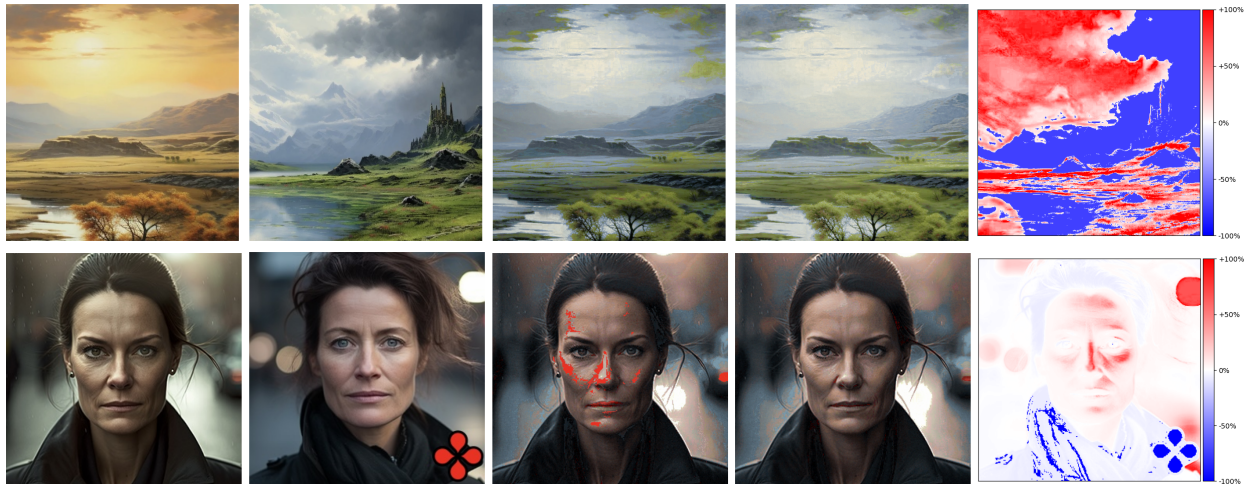


Figure 3: **Color transfer** between a source and a target image (*first and second columns*). We compare SOT gradient flows operated in the color space (*third column*) and the same procedure with a reweighing of the distributions by USOT (*fourth column*). The last column shows a percentage of mass change given by USOT, *i.e.*, $\frac{(\pi_2^* - \beta)}{\beta}$, where *red* indicates mass creation and *blue* mass destruction.

of this experiment is to show that by discarding possible outliers using a well chosen parameter ρ , USOT is able to outperform SOT and SUOT on this task. Since a word typically appears several times in a document, the measures are not uniform and sliced partial OT (Bonneel & Coeurjolly, 2019; Bai et al., 2023) cannot be used in this setting. We detail in Appendix C additional experiments without normalizing histograms to compare with (Bai et al., 2023) (*i.e.*, $\sum_{i=1}^{n_k} w_i^k$ is the sentence length). We consider the BBCSport dataset (Kusner et al., 2015), a standard benchmark with small documents for which OT can be used effectively, and the Goodreads dataset (Maharjan et al., 2017) on two tasks (genre and likability predictions), a dataset with large-scale documents for which the computational burden of performing OT and UOT is substantial. We report on Table 1 the accuracy and average runtimes of OT, UOT computed with the majorization minimization algorithm (Chapel et al., 2021) or approximated with the Sinkhorn algorithm (SinkhUOT) (Pham et al., 2020), as well as SOT, SUOT and USOT. All the benchmark methods are computed using the Python OT library (Flamary et al., 2021) on a Nvidia Tesla V100 GPU. For sliced methods, we average over 3 computations of the loss matrix and report the standard deviation in Table 1. The number of neighbors was selected via cross validation. The results for UOT, SinkhUOT, SUOT and USOT are reported for ρ yielding the best accuracy among a grid (see Appendix C.1 for more details), and we display an ablation of this parameter on the BBCSport dataset in Figure 2. We also add on Table 1 the results obtained with a cross validation (CV) on ρ for USOT and SUOT. Our findings demonstrate that our approaches surpass SOT in performance, incurring only a minor computational overhead. Moreover, our methods closely rival the performance of OT, while being 40 times faster on large-scale datasets. This highlight their practical significance, particularly when OT is computationally unfeasible.

5.2 Color transfer

Color transfer is a long-standing problem in OT, which dates back to the seminal work of Rabin et al. (2010). It consists in aligning the color distributions of two images. While previous works, *e.g.* (Ferradans et al., 2013; Bonneel et al., 2016), considered color palettes to deal with the complexity of OT, we illustrate the scalability of our methods by considering here the full distributions of pixels within images, in a way similar to (Bonnel & Coeurjolly, 2019). We express the color transfer as a gradient flow, where every pixel is a sample in the 3D RGB color space. Formally, let $\alpha(t) = \frac{1}{N} \sum_{i=1}^N \delta_{x_i(t)}$, $\beta = \frac{1}{M} \sum_{j=1}^M \delta_{y_j}$, where α (resp. β) represents the color distribution of the source (resp. target) image. The SOT gradient flow performing color transfer consists in iterating the following scheme: $X(t+1) = X(t) - \gamma \nabla_X \text{SOT}(\alpha(t), \beta)$, where $\alpha(t)$ is the color distribution of the source image at iterations t , supported by pixels from $X(t)$. One of the major

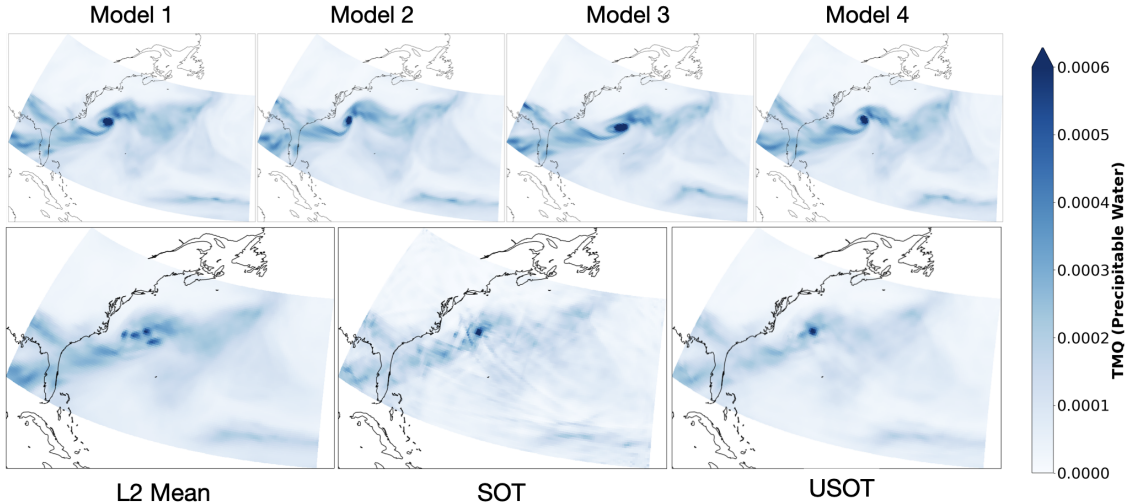


Figure 4: **Barycenter of geophysical data.** (*First row*) Simulated output of 4 different climate models depicting different scenarios for the evolution of a tropical cyclone (*Second row*) Results of different aggregation strategies.

problem is color bleeding: potential color artefacts are likely to appear since the object proportions between images are likely to differ. We propose to correct this issue in a two steps procedure, relying on USOT: *i*) we first obtain optimal marginals (π_1^*, π_2^*) by solving (9) and using the Norm routine, then *ii*) we solve the classical SOT gradient flow with the measures (π_1^*, π_2^*) .

We present results on 300×300 images produced by the generative model Midjourney (Mid, 2023). The images correspond to two landscapes and photo-realistic portraits of two women (see first and second column of Figure 3). For every results, we iterate the gradient flow for 100 iterations, and the learning rate γ is set to 10^{-2} . The computation of the final result is produced in less than one minute with a commodity GPU, while it is out of reach for OT or UOT solvers on this distributions size (90K pixels). The third column shows the color transfer using SOT. It reveals instances of color bleeding: green clouds appear in the sky of landscape scenes, and a red superimposed logo results in unwanted red pixels in the portrait. Contrasting, the fourth column displays the outcomes achieved through our approach, effectively addressing the color bleeding issue. We chose $\rho_1 = 10^4$ and $\rho_2 = 0.02$ for this task and we observe the resulting change in mass distribution on the target image in the last column of Figure 3. It provides insightful indications of the discarded colors (the darkened landscape in the first case and the removal of the logo in the second), which is only possible when all the pixels are considered in the transfer.

5.3 Barycenter of geophysical data

OT barycenters are an important topic of interest (Le et al., 2021) for their ability to capture mass changes and spatial deformations over several reference measures. In order to compute barycenters under the USOT geometry on a fixed grid, we employ a mirror-descent strategy similar to (Cuturi & Doucet, 2014a, Algorithm (1)) and described more in depth in Appendix C. We compute unbalanced sliced OT barycenter for climate model data. Ensembles of multiple models are commonly employed to reduce biases and evaluate uncertainties in climate projections (Sanderson et al., 2015; Thao et al., 2022). The commonly used Multi-Model Mean approach assumes models are centered around true values and averages the ensemble with equal or varying weights. However, spatial averaging may fail in capturing specific characteristics of the physical system at stake, and we propose to use USOT barycenter instead. We use the ClimateNet dataset (Prabhat et al., 2021), and more specifically the TMQ (precipitable water) indicator. The ClimateNet dataset is a human-expert-labeled curated dataset which captures tropical cyclones (TCs), among other things. To simulate the output of several climate models, we take a specific instant (first date of 2011) and apply the elastic deformation from TorchVision (Paszke et al., 2019) in an area close to the eastern part of the U.S.A. As a result, we obtain 4 different TCs, as shown in the first row of Figure 4. The classical L2 spatial mean is

displayed on the second row of Figure 4 and reveals 4 different TCs centers/modes, which is undesirable. As the total TMQ mass in the considered zone varies between the different models, a direct application of SOT is impossible, or requires a normalization of the mass that has undesired effect as can be seen on the second row. Finally, we show the result of the USOT barycenter with $\rho_1 = 1e1$ (related to the data) and $\rho_2 = 1e4$ (related to the barycenter). This barycenter has only one apparent mode, which is the expected behaviour. The considered measures have a size of 100×200 , and we run the barycenter algorithm for 500 iterations (with $K = 64$ projections), which takes 3 minutes on a commodity GPU. UOT barycenters for this size of problems are intractable, and to the best of our knowledge, our experiment is the very first instance where unbalanced OT barycenters can be computed on such a large scale.

6 Conclusion

We proposed two losses merging unbalanced and sliced OT, with theoretical guarantees and an efficient and modular Frank-Wolfe algorithm. We illustrate the performance improvement over SOT on various experiments, and described novel applications of unbalanced OT barycenters of positive measures, with a new case study on geophysical data. These novel results and algorithms pave the way to numerous new applications of sliced variants of OT, and we believe that our contributions will motivate practitioners to further explore their use in general ML applications, without the cumbersome task of pre-processing probability measures.

Acknowledgments

We thank the anonymous reviewers for their valuable comments. KF was supported by NSERC Discovery grant (RGPIN-2019-06512) and a Samsung grant. CB was supported by project DynaLearn from Labex CominLabs and Région Bretagne ARED DLearnMe, and by the ANR PEPR PDE-AI. NC was supported by the ANR AI Chair OTTOPIA ANR-20-CHIA-0030

References

- Midjourney, 2023. URL "<https://docs.midjourney.com/legacy/en>". (Cited on p. 13)
- Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media, 2005. (Cited on p. 8)
- Yikun Bai, Bernhard Schmitzer, Matthew Thorpe, and Soheil Kolouri. Sliced optimal partial transport. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13681–13690, 2023. (Cited on p. 2, 4, 5, 8, 9, 12, 39)
- Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2):423–443, 2018. (Cited on p. 1)
- Erhan Bayraktar and Gaoyue Guo. Strong equivalence between metrics of Wasserstein type. *Electronic Communications in Probability*, 26(none):1 – 13, 2021. doi: 10.1214/21-ECP383. (Cited on p. 2, 8)
- Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003. (Cited on p. 41)
- Vladimir Igorevich Bogachev and Maria Aparecida Soares Ruas. *Measure theory*, volume 1. Springer, 2007. (Cited on p. 23)
- Clément Bonet, Benoît Malézieux, Alain Rakotomamonjy, Lucas Drumetz, Thomas Moreau, Matthieu Kowalski, and Nicolas Courty. Sliced-wasserstein on symmetric positive definite matrices for m/eeg signals. In *Proceedings of the 40th International Conference on Machine Learning*, 2023a. (Cited on p. 2)
- Clément Bonet, Nicolas Courty, François Septier, and Lucas Drumetz. Efficient gradient flows in sliced-wasserstein space. *Transactions on Machine Learning Research*, 2022. (Cited on p. 8)

- Clément Bonet, Paul Berg, Nicolas Courty, François Septier, Lucas Drumetz, and Minh-Tan Pham. Spherical sliced-wasserstein. In *The Eleventh International Conference on Learning Representations*, 2023b. (Cited on p. 2)
- Clément Bonet, Laetitia Chapel, Lucas Drumetz, and Nicolas Courty. Hyperbolic sliced-wasserstein via geodesic and horospherical projections. In *Proceedings of 2nd Annual Workshop on Topology, Algebra, and Geometry in Machine Learning (TAG-ML)*, pp. 334–370. PMLR, 2023c. (Cited on p. 2, 11, 41)
- Nicolas Bonneel and David Coeurjolly. Spot: sliced partial optimal transport. *ACM Transactions on Graphics (TOG)*, 38(4):1–13, 2019. (Cited on p. 2, 4, 5, 8, 12)
- Nicolas Bonneel, Julien Rabin, Gabriel Peyré, and Hanspeter Pfister. Sliced and radon wasserstein barycenters of measures. *Journal of Mathematical Imaging and Vision*, 51:22–45, 2015. (Cited on p. 2)
- Nicolas Bonneel, Gabriel Peyré, and Marco Cuturi. Wasserstein barycentric coordinates: histogram regression using optimal transport. *ACM Trans. Graph.*, 35(4), jul 2016. (Cited on p. 12)
- Nicolas Bonnotte. *Unidimensional and evolution methods for optimal transportation*. PhD thesis, Université Paris Sud-Paris XI; Scuola normale superiore (Pise, Italie), 2013. (Cited on p. 8, 25, 26)
- Jules Candau-Tilh. Wasserstein and sliced-wasserstein distances. Master’s thesis, Université Pierre et Marie Curie, 2020. (Cited on p. 8)
- Laetitia Chapel, Rémi Flamary, Haoran Wu, Cédric Févotte, and Gilles Gasso. Unbalanced optimal transport through non-negative penalized linear regression. *Advances in Neural Information Processing Systems*, 34: 23270–23282, 2021. (Cited on p. 12, 37)
- Lenaïc Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. *Advances in neural information processing systems*, 31, 2018. (Cited on p. 2, 8)
- Lenaïc Chizat, Gabriel Peyré, Bernhard Schmitzer, and François-Xavier Vialard. An interpolating distance between optimal transport and fisher-rao metrics. *Foundations of Computational Mathematics*, 18:1–44, 2018a. (Cited on p. 3)
- Lenaïc Chizat, Gabriel Peyré, Bernhard Schmitzer, and François-Xavier Vialard. Unbalanced optimal transport: Dynamic and kantorovich formulations. *Journal of Functional Analysis*, 274(11):3090–3123, 2018b. (Cited on p. 2)
- Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems*, 26, 2013. (Cited on p. 2)
- Marco Cuturi and Arnaud Doucet. Fast computation of wasserstein barycenters. In Eric P. Xing and Tony Jebara (eds.), *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pp. 685–693, Beijing, China, 22–24 Jun 2014a. PMLR. (Cited on p. 13)
- Marco Cuturi and Arnaud Doucet. Fast computation of wasserstein barycenters. In *International conference on machine learning*, pp. 685–693. PMLR, 2014b. (Cited on p. 41)
- Pinar Demetci, Rebecca Santorella, Manav Chakravarthy, Bjorn Sandstede, and Ritambhara Singh. Scotv2: Single-cell multiomic alignment with disproportionate cell-type representation. *Journal of Computational Biology*, 29(11):1213–1228, 2022. (Cited on p. 2)
- Ishan Deshpande, Yuan-Ting Hu, Ruoyu Sun, Ayis Pyrros, Nasir Siddiqui, Sanmi Koyejo, Zhizhen Zhao, David Forsyth, and Alexander G Schwing. Max-sliced wasserstein distance and its use for gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10648–10656, 2019. (Cited on p. 2, 10)

- Richard Mansfield Dudley. The speed of mean glivenko-cantelli convergence. *The Annals of Mathematical Statistics*, 40(1):40–50, 1969. (Cited on p. 2)
- Kilian Fatras, Younes Zine, Rémi Flamary, Remi Gribonval, and Nicolas Courty. Learning with minibatch wasserstein : asymptotic and gradient properties. In Silvia Chiappa and Roberto Calandra (eds.), *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pp. 2131–2141, Online, 26–28 Aug 2020. PMLR. (Cited on p. 2)
- Kilian Fatras, Thibault Sejourne, Rémi Flamary, and Nicolas Courty. Unbalanced minibatch optimal transport; applications to domain adaptation. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 3186–3197. PMLR, 18–24 Jul 2021. (Cited on p. 1, 2)
- Sira Ferradans, Nicolas Papadakis, Julien Rabin, Gabriel Peyré, and Jean-François Aujol. Regularized discrete optimal transport. In *Scale Space and Variational Methods in Computer Vision*, pp. 428–439, 2013. (Cited on p. 12)
- Rémi Flamary, Nicolas Courty, Alexandre Gramfort, Mokhtar Z Alaya, Aurélie Boisbunon, Stanislas Chambon, Laetitia Chapel, Adrien Corenflos, Kilian Fatras, Nemo Fournier, et al. Pot: Python optimal transport. *The Journal of Machine Learning Research*, 22(1):3571–3578, 2021. (Cited on p. 12)
- Aude Genevay, Lénaïc Chizat, Francis Bach, Marco Cuturi, and Gabriel Peyré. Sample complexity of sinkhorn divergences. In *The 22nd international conference on artificial intelligence and statistics*, pp. 1574–1583. PMLR, 2019. (Cited on p. 30)
- Ziv Goldfeld and Kristjan Greenewald. Sliced mutual information: A scalable measure of statistical dependence. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 17567–17578. Curran Associates, Inc., 2021. (Cited on p. 2)
- Elad Hazan and Haipeng Luo. Variance-reduced and projection-free stochastic optimization. In *International Conference on Machine Learning*, pp. 1263–1271. PMLR, 2016. (Cited on p. 11)
- Soheil Kolouri, Kimia Nadjahi, Umut Simsekli, Roland Badeau, and Gustavo Rohde. Generalized sliced wasserstein distances. *Advances in neural information processing systems*, 32, 2019. (Cited on p. 2)
- Stanislav Kondratyev, Léonard Monsaingeon, and Dmitry Vorotnikov. A fitness-driven cross-diffusion system from population dynamics as a gradient flow. *Journal of Differential Equations*, 261(5):2784–2808, 2016. (Cited on p. 2)
- Matt Kusner, Yu Sun, Nicholas Kolkin, and Kilian Weinberger. From word embeddings to document distances. In *International conference on machine learning*, pp. 957–966. PMLR, 2015. (Cited on p. 11, 12, 37)
- Simon Lacoste-Julien and Martin Jaggi. On the global linear convergence of frank-wolfe optimization variants. *Advances in neural information processing systems*, 28, 2015. (Cited on p. 10)
- Khang Le, Huy Nguyen, Quang M Nguyen, Tung Pham, Hung Bui, and Nhat Ho. On robust optimal transport: Computational complexity and barycenter computation. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 21947–21959, 2021. (Cited on p. 13)
- Tam Le and Truyen Nguyen. Entropy partial transport with tree metrics: Theory and practice. In *International Conference on Artificial Intelligence and Statistics*, pp. 3835–3843. PMLR, 2021. (Cited on p. 2)
- Matthias Liero, Alexander Mielke, and Giuseppe Savaré. Optimal entropy-transport problems and a new hellinger–kantorovich distance between positive measures. *Inventiones mathematicae*, 211(3):969–1117, 2018. (Cited on p. 2, 3, 4, 5, 6, 8, 19, 21, 29, 30, 36)

- Suraj Maharjan, John Arevalo, Manuel Montes, Fabio A González, and Tamar Solorio. A multi-task approach to predict likability of books. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pp. 1217–1227, 2017. (Cited on p. 12, 37)
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26, 2013. (Cited on p. 11, 37)
- Kimia Nadjahi, Alain Durmus, Lénaïc Chizat, Soheil Kolouri, Shahin Shahrampour, and Umut Simsekli. Statistical and topological properties of sliced probability divergences. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 20802–20812. Curran Associates, Inc., 2020a. (Cited on p. 5, 6)
- Kimia Nadjahi, Alain Durmus, Lénaïc Chizat, Soheil Kolouri, Shahin Shahrampour, and Umut Simsekli. Statistical and topological properties of sliced probability divergences. *Advances in Neural Information Processing Systems*, 33:20802–20812, 2020b. (Cited on p. 2, 6, 24, 32)
- Yoshihiro Nagano, Shoichiro Yamaguchi, Yasuhiro Fujita, and Masanori Koyama. A wrapped normal distribution on hyperbolic space for gradient-based learning. In *International Conference on Machine Learning*, pp. 4693–4702. PMLR, 2019. (Cited on p. 41)
- Khai Nguyen, Nhat Ho, Tung Pham, and Hung Bui. Distributional sliced-wasserstein and applications to generative modeling. *arXiv preprint arXiv:2002.07367*, 2020. (Cited on p. 2, 10)
- Khai Nguyen, Tongzheng Ren, Huy Nguyen, Litu Rout, Tan Minh Nguyen, and Nhat Ho. Hierarchical sliced wasserstein distance. In *The Eleventh International Conference on Learning Representations*, 2023. (Cited on p. 2)
- Ruben Ohana, Kimia Nadjahi, Alain Rakotomamonjy, and Liva Ralaivola. Shedding a pac-bayesian light on adaptive sliced-wasserstein distances. In *Proceedings of the 40th International Conference on Machine Learning*, 2023. (Cited on p. 2)
- Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. Thumbs up? sentiment classification using machine learning techniques. In *Proceedings of EMNLP*, pp. 79–86, 2002. (Cited on p. 37)
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. (Cited on p. 10, 13)
- Ofir Pele and Michael Werman. Fast and robust earth mover’s distances. In *2009 IEEE 12th international conference on computer vision*, pp. 460–467. IEEE, 2009. (Cited on p. 2)
- Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport: With applications to data science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019. (Cited on p. 2, 4)
- Khiem Pham, Khang Le, Nhat Ho, Tung Pham, and Hung Bui. On unbalanced optimal transport: An analysis of Sinkhorn algorithm. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 7673–7682. PMLR, 13–18 Jul 2020. (Cited on p. 2, 4, 12)
- Benedetto Piccoli and Francesco Rossi. Generalized wasserstein distance and its application to transport equations with source. *Archive for Rational Mechanics and Analysis*, 211:335–358, 2014. (Cited on p. 30)
- Prabhat, K. Kashinath, M. Mudigonda, S. Kim, L. Kapp-Schwoerer, A. Graubner, E. Karaismailoglu, L. von Kleist, T. Kurth, A. Greiner, A. Mahesh, K. Yang, C. Lewis, J. Chen, A. Lou, S. Chandran, B. Toms, W. Chapman, K. Dagon, C. A. Shields, T. O’Brien, M. Wehner, and W. Collins. Climateset: an expert-labeled open dataset and deep learning architecture for enabling high-precision analyses of extreme weather. *Geoscientific Model Development*, 14(1):107–124, 2021. doi: 10.5194/gmd-14-107-2021. (Cited on p. 13)

- Julien Rabin, Julie Delon, and Yann Gousseau. Regularization of transportation maps for color and contrast transfer. In *International Conference on Image Processing*, pp. 1933–1936, 2010. (Cited on p. 12)
- Julien Rabin, Gabriel Peyré, Julie Delon, and Marc 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*, pp. 435–446. Springer, 2012. (Cited on p. 2, 4)
- Grant Rotskoff, Samy Jelassi, Joan Bruna, and Eric Vanden-Eijnden. Global convergence of neuron birth-death dynamics. *arXiv preprint arXiv:1902.01843*, 2019. (Cited on p. 2, 8)
- Benjamin M Sanderson, Reto Knutti, and Peter Caldwell. A representative democracy to reduce interdependency in a multimodel ensemble. *Journal of Climate*, 28(13):5171–5194, 2015. (Cited on p. 13)
- Filippo Santambrogio. Optimal transport for applied mathematicians. *Birkhäuser, NY*, 55(58-63):94, 2015. (Cited on p. 19, 23, 29, 31, 34, 35)
- Ryoma Sato, Makoto Yamada, and Hisashi Kashima. Fast unbalanced optimal transport on a tree. *Advances in neural information processing systems*, 33:19039–19051, 2020. (Cited on p. 2)
- Geoffrey Schiebinger, Jian Shu, Marcin Tabaka, Brian Cleary, Vidya Subramanian, Aryeh Solomon, Joshua Gould, Siyan Liu, Stacie Lin, Peter Berube, et al. Optimal-transport analysis of single-cell gene expression identifies developmental trajectories in reprogramming. *Cell*, 176(4):928–943, 2019. (Cited on p. 1, 2)
- Thibault Séjourné, Jean Feydy, François-Xavier Vialard, Alain Trounev, and Gabriel Peyré. Sinkhorn divergences for unbalanced optimal transport. *arXiv preprint arXiv:1910.12958*, 2019. (Cited on p. 21)
- Thibault Séjourné, Gabriel Peyré, and François-Xavier Vialard. Unbalanced optimal transport, from theory to numerics. *arXiv preprint arXiv:2211.08775*, 2022a. (Cited on p. 2, 3, 5)
- Thibault Séjourné, François-Xavier Vialard, and Gabriel Peyré. Faster unbalanced optimal transport: Translation invariant sinkhorn and 1-d frank-wolfe. In *International Conference on Artificial Intelligence and Statistics*, pp. 4995–5021. PMLR, 2022b. (Cited on p. 2, 9, 23, 29, 34)
- Stephen Simons. *Minimax and monotonicity*. Springer, 2006. (Cited on p. 21)
- Soulihanh Thao, Mats Garvik, Gregoire Mariethoz, and Mathieu Vrac. Combining global climate models using graph cuts. *Climate Dynamics*, 59:2345–2361, 2022. (Cited on p. 13)
- Adrien Vacher and François-Xavier Vialard. Semi-dual unbalanced quadratic optimal transport: fast statistical rates and convergent algorithm. In *International Conference on Machine Learning*, pp. 34734–34758. PMLR, 2023. (Cited on p. 6)
- Jiaqi Xi and Jonathan Niles-Weed. Distributional convergence of the sliced wasserstein process. *arXiv preprint arXiv:2206.00156*, 2022. (Cited on p. 23)

A Postponed proofs for Section 3

A.1 Existence of minimizers: Proof of Proposition 3.2 and Proposition 3.7

We provide the formal statement and detailed proof on the existence of a solution for both SUOT and USOT, as mentioned in Section 3.

Proposition A.1. (Existence of minimizers) *Assume that C_1 is lower-semicontinuous and that either (i) $\varphi'_{1,\infty} = \varphi'_{2,\infty} = +\infty$, or (ii) C_1 has compact sublevels on $\mathbb{R} \times \mathbb{R}$ and $\varphi'_{1,\infty} + \varphi'_{2,\infty} + \inf C_1 > 0$. Then the solution of SUOT(α, β) and USOT(α, β) exist, i.e., the infimum in (6) and (9) is attained. More precisely, there exists (π_1, π_2) which attains the infimum for USOT(α, β) (see Equation 9). Concerning SUOT(α, β), there exists for any $\theta \in \text{supp}(\sigma)$ a plan π_θ attaining the infimum in UOT($\theta_\#^* \alpha, \theta_\#^* \beta$) (see Equation 2).*

Proof. We leverage (Liero et al., 2018, Theorem 3.3) to prove this proposition. In the setting of SUOT, if such assumptions (i) or (ii) are satisfied for (α, β) , then they also hold for $(\theta_\#^* \alpha, \theta_\#^* \beta)$ for any $\theta \in \mathbb{S}^{d-1}$. Hence, UOT($\theta_\#^* \alpha, \theta_\#^* \beta$) admits a solution π^θ .

Concerning USOT, note that one necessarily has $m(\pi_1) = m(\pi_2)$, otherwise SOT(π_1, π_2) = $+\infty$. From (Liero et al., 2018, Equation (3.10)), for any admissible (π_1, π_2, π) , one has

$$\text{USOT}(\alpha, \beta) \geq m(\pi) \inf C_1 + m(\alpha) \varphi_1\left(\frac{m(\pi)}{m(\alpha)}\right) + m(\beta) \varphi_2\left(\frac{m(\pi)}{m(\beta)}\right).$$

In both settings the above bounds implies coercivity of the functional of USOT w.r.t. the masses of the measures (π_1, π_2, π) . Thus there exists $M > 0$ such that $m(\pi_1) = m(\pi_2) = m(\pi) < M$, otherwise USOT(α, β) = $+\infty$. By the Banach-Alaoglu theorem, the set of bounded measures (π_1, π_2) is compact, and the set of plans π with such marginals is also compact because \mathbb{R}^d is Polish and C_1 is lower-semicontinuous (Santambrogio, 2015, Theorem 1.7). Because the functional of USOT is lower-semicontinuous in (π_1, π_2, π) and we can restrict optimization over a compact set, we have existence of minimizers for USOT by standard proofs of calculus of variations. \square

A.2 Strong duality: Proof of Theorem 3.5 and Theorem 3.9

Note that the result for SUOT (Theorem 3.5) is proved in Lemma A.4. Thus we focus on the proof of duality for USOT.

Proof of Theorem 3.9. We start from the definition of USOT, reformulate it to apply the strong duality result of Proposition A.2 and obtain our reformulation. We first have that

$$\begin{aligned} \text{USOT}(\alpha, \beta) &= \inf_{(\pi_1, \pi_2) \in \mathcal{M}_+(\mathbb{R}^d)^2} \{ \text{SOT}(\pi_1, \pi_2) + D_{\varphi_1}(\pi_1 | \alpha) + D_{\varphi_2}(\pi_2 | \beta) \}, \\ &= \inf_{(\pi_1, \pi_2) \in \mathcal{M}_+(\mathbb{R}^d)^2} \left\{ \int_{\mathbb{S}^{d-1}} \left[\sup_{f_\theta \oplus g_\theta \leq C_1} \int f_\theta d(\theta_\#^* \pi_1) + \int g_\theta d(\theta_\#^* \pi_2) \right] d\hat{\sigma}_K(\theta) \right. \\ &\quad \left. + \sup_{\tilde{f} \in \mathcal{E}(\mathbb{R}^d)} \int \varphi_1^\circ(\tilde{f}(x)) d\alpha(x) - \int \tilde{f}(x) d\pi_1(x) \right. \\ &\quad \left. + \sup_{\tilde{g} \in \mathcal{E}(\mathbb{R}^d)} \int \varphi_2^\circ(\tilde{g}(y)) d\beta(y) - \int \tilde{g}(y) d\pi_2(y) \right\}, \\ &= \inf_{(\pi_1, \pi_2) \in \mathcal{M}_+(\mathbb{R}^d)^2} \left\{ \sup_{f_\theta \oplus g_\theta \leq C_1} \int_{\mathbb{S}^{d-1}} \left[\int f_\theta d(\theta_\#^* \pi_1) + \int g_\theta d(\theta_\#^* \pi_2) \right] d\hat{\sigma}_K(\theta) \right. \\ &\quad \left. + \sup_{\tilde{f} \in \mathcal{E}(\mathbb{R}^d)} \int \varphi_1^\circ(\tilde{f}(x)) d\alpha(x) - \int \tilde{f}(x) d\pi_1(x) \right. \\ &\quad \left. + \sup_{\tilde{g} \in \mathcal{E}(\mathbb{R}^d)} \int \varphi_2^\circ(\tilde{g}(y)) d\beta(y) - \int \tilde{g}(y) d\pi_2(y) \right\}, \end{aligned}$$

where $\mathcal{E}(\mathbb{R}^d)$ denotes a set of lower-semicontinuous functions, and the last equality holds thanks to Lemma A.3.

We focus now on verifying that Proposition A.2 holds, so that we can swap the infimum and the supremum. Define the functional

$$\begin{aligned} \mathcal{L}((\pi_1, \pi_2), ((f_\theta)_\theta, (g_\theta)_\theta, \tilde{f}, \tilde{g})) &\triangleq \int_{\mathbb{S}^{d-1}} \left[\int f_\theta d(\theta_\#^* \pi_1) + \int g_\theta d(\theta_\#^* \pi_2) \right] d\hat{\sigma}_K(\theta) \\ &\quad + \int \varphi_1^\circ(\tilde{f}(x)) d\alpha(x) - \int \tilde{f}(x) d\pi_1(x) \\ &\quad + \int \varphi_2^\circ(\tilde{g}(y)) d\beta(y) - \int \tilde{g}(y) d\pi_2(y). \end{aligned}$$

One has that,

- For any $((f_\theta)_\theta, (g_\theta)_\theta, \tilde{f}, \tilde{g})$, \mathcal{L} is linear (thus convex) and lower-semicontinuous.
- For any (π_1, π_2) , \mathcal{L} is concave in $((f_\theta)_\theta, (g_\theta)_\theta, \tilde{f}, \tilde{g})$ because φ_i° is concave and thus \mathcal{L} is a sum of linear or concave functions.

Furthermore, since we assumed that $0 \in \text{dom}(\varphi)$, then

$$\sup_{((f_\theta)_\theta, (g_\theta)_\theta, \tilde{f}, \tilde{g})} \inf_{(\pi_1, \pi_2) \in \mathcal{M}_+(\mathbb{R}^d)^2} \mathcal{L} \leq \text{USOT}(\alpha, \beta) \leq \varphi_1(0)m(\alpha) + \varphi_2(0)m(\beta),$$

because the marginals $(\pi_1, \pi_2) = (0, 0)$ are admissible and suboptimal. If we consider instead that $(m(\alpha), m(\beta)) \in \text{dom}(\varphi)$, then we take the marginals $\pi_1 = \alpha/m(\alpha)$ and $\pi_2 = \beta/m(\beta)$, which yields an upper-bound by $m(\alpha)\varphi_1(\frac{1}{m(\alpha)}) + m(\beta)\varphi_2(\frac{1}{m(\beta)})$. Then we consider an anchor dual point $b^* = ((f_\theta)_\theta, (g_\theta)_\theta, \tilde{f}, \tilde{g})$ to bound \mathcal{L} over a compact set. We take $f_\theta = 0, g_\theta = 0$, which are always admissible since we take $C_1(x, y) \geq 0$. Then, since we assume there exists $p_i \leq 0$ in $\text{dom}(\varphi_i^*)$, we take $\tilde{f} = p_1$ and $\tilde{g} = p_2$. For these potentials one has:

$$\mathcal{L}((\pi_1, \pi_2), b^*) = \varphi_1^\circ(p_1)m(\alpha) - p_1m(\pi_1) + \varphi_2^\circ(p_2)m(\beta) - p_2m(\pi_2).$$

Note that the functional at this point only depends on the masses of the marginals (π_1, π_2) . Since $(p_1, p_2) \geq 0$ the set of (π_1, π_2) such that $\mathcal{L}((\pi_1, \pi_2), b^*) \leq \varphi_1(0)m(\alpha) + \varphi_2(0)m(\beta)$ is non-empty (at least in a neighbourhood of $(\pi_1, \pi_2) = (0, 0)$), and that $(m(\pi_1), m(\pi_2))$ are uniformly bounded by some constant $M > 0$. By the Banach-Alaoglu theorem, such set of measures is compact for the weak* topology.

Therefore, Proposition A.2 holds and we have strong duality, *i.e.*,

$$\text{USOT}(\alpha, \beta) = \sup_{\left\{ \begin{array}{l} f_\theta \oplus g_\theta \leq C_1 \\ (\tilde{f}, \tilde{g}) \in \mathcal{E}(\mathbb{R}^d) \end{array} \right\}} \inf_{(\pi_1, \pi_2) \in \mathcal{M}_+(\mathbb{R}^d)^2} \mathcal{L}((\pi_1, \pi_2), ((f_\theta)_\theta, (g_\theta)_\theta, \tilde{f}, \tilde{g})).$$

To achieve the proof, note that taking the infimum in (π_1, π_2) (for fixed dual variables) reads

$$\begin{aligned} \inf_{\pi_1, \pi_2 \geq 0} \int &\left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\hat{\sigma}_K(\theta) \right) d\pi_1(x) - \int \tilde{f}(x) d\pi_1(x) \\ &+ \int \left(\int_{\mathbb{S}^{d-1}} g_\theta(\theta^*(y)) d\hat{\sigma}_K(\theta) \right) d\pi_2(y) - \int \tilde{g}(y) d\pi_2(y). \end{aligned}$$

Note that we applied Fubini's theorem here, which holds here because all measures have compact support, thus all quantities are finite. It allows to rephrase the minimization over $\pi_1, \pi_2 \geq 0$ as the following constraint

$$\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\hat{\sigma}_K(\theta) \geq \tilde{f}(x), \quad \int_{\mathbb{S}^{d-1}} g_\theta(\theta^*(y)) d\hat{\sigma}_K(\theta) \geq \tilde{g}(y),$$

otherwise the infimum is $-\infty$. However, the function φ° is non-decreasing (see (Séjourné et al., 2019, Proposition 2)). Thus the maximization in (\tilde{f}, \tilde{g}) is optimal when the above inequality is actually an equality, *i.e.*,

$$\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\hat{\sigma}_K(\theta) = \tilde{f}(x), \quad \int_{\mathbb{S}^{d-1}} g_\theta(\theta^*(y)) d\hat{\sigma}_K(\theta) = \tilde{g}(y).$$

Plugging the above relation in the functional \mathcal{L} yields the desired result on the dual of USOT and ends the proof. \square

We mention a strong duality result which is very general and which we use in the proof of Theorem 3.9. This result is taken from (Liero et al., 2018, Theorem 2.4) which itself takes it from (Simons, 2006).

Proposition A.2. (Liero et al., 2018, Theorem 2.4) *Consider two sets A and B be nonempty convex sets of some vector spaces. Assume A is endowed with a Hausdorff topology. Let $L : A \times B \rightarrow \mathbb{R}$ be a function such that*

1. $a \mapsto L(a, b)$ is convex and lower-semicontinuous on A , for every $b \in B$
2. $b \mapsto L(a, b)$ is concave on B , for every $a \in A$.

If there exists $b_\star \in B$ and $\kappa > \sup_{b \in B} \inf_{a \in A} L(a, b)$ such that the set $\{a \in A, L(a, b_\star) < \kappa\}$ is compact in A , then

$$\inf_{a \in A} \sup_{b \in B} L(a, b) = \sup_{b \in B} \inf_{a \in A} L(a, b)$$

We also consider the following to swap the supremum in the integral which defines sliced-UOT (and in particular sliced-OT). In what follows we note sliced potentials as functions $f_\theta(z)$ with $(\theta, z) \in \mathbb{S}^{d-1} \times \mathbb{R}$, such that

$$\text{SUOT}(\alpha, \beta) = \int_{\mathbb{S}^{d-1}} \left[\sup_{f_\theta \oplus g_\theta \leq C_1} \int \varphi^\circ \circ f_\theta d(\theta_\#^* \alpha) + \int \varphi^\circ \circ g_\theta d(\theta_\#^* \beta) \right] d\hat{\sigma}_K(\theta).$$

Note that with the above definition, $z \mapsto f_\theta(z)$ is continuous for any θ , but $\theta \mapsto f_\theta(z)$ is only $\hat{\sigma}_K$ -measurable.

Lemma A.3. *Consider two sets X and Y , a measure σ such that $\sigma(X) < +\infty$. Assume Y is compact. Consider a function $\mathcal{F} : X \times Y \rightarrow \mathbb{R}$. Assume there exists a sequence (y_n) in Y such that $\mathcal{F}(\cdot, y_n) \rightarrow \sup_{y \in Y} \mathcal{F}(\cdot, y)$ uniformly. Then one has*

$$\sup_{y \in Y} \int_X \mathcal{F}(x, y) d\sigma(x) = \int_X \sup_{y \in Y} \mathcal{F}(x, y) d\sigma(x).$$

Proof. Define $\mathcal{G}(x) = \sup_{y \in Y} \mathcal{F}(x, y)$ and $\mathcal{H}(x, y) \triangleq \mathcal{G}(x) - \mathcal{F}(x, y)$. One has $\mathcal{H} \geq 0$ by definition, and the desired equality can be rewritten as

$$\begin{aligned} \sup_{y \in Y} \int_X \mathcal{F}(x, y) d\sigma(x) &= \int_X \sup_{y \in Y} \mathcal{F}(x, y) d\sigma(x) \\ \Leftrightarrow \inf_{y \in Y} \int_X \mathcal{H}(x, y) d\sigma(x) &= 0. \end{aligned}$$

Since the integral involving \mathcal{H} is non-negative, the infimum is zero if and only if we have a sequence (y_n) such that $\int_X \mathcal{H}(\cdot, y_n) d\sigma \rightarrow 0$. By assumption, one has $\mathcal{F}(\cdot, y_n) \rightarrow \sup_{y \in Y} \mathcal{F}(\cdot, y)$ uniformly, *i.e.*, $\|\mathcal{H}(\cdot, y_n)\|_\infty \rightarrow 0$. This implies thanks to Holder's inequality that

$$0 \leq \int_X \mathcal{H}(\cdot, y_n) d\sigma \leq \sigma(X) \|\mathcal{H}(\cdot, y_n)\|_\infty$$

Thus by assumption one has $\int_X \mathcal{F}(\cdot, y_n) d\sigma \rightarrow \int_X \mathcal{G} d\sigma$, which indeed means that we have the desired permutation between supremum and integral. \square

Lemma A.4. *Let $p \in [1, +\infty)$ and assume that $C_1(x, y) = |x - y|^p$. Consider two positive measures (α, β) with compact support. Assume that the measure $\hat{\sigma}_K$ is discrete, i.e., $\hat{\sigma}_K = \frac{1}{K} \sum_{i=1}^K \delta_{\theta_i}$ with $\theta_i \in \mathbb{S}^{d-1}$, $i = 1, \dots, n$. Then, one can swap the integral over the sphere and the supremum in the dual formulation of SUOT, such that*

$$\text{SUOT}(\alpha, \beta) = \sup_{f_\theta \oplus g_\theta \leq C_1} \int_{\mathbb{S}^{d-1}} \left[\int \varphi^\circ \circ f_\theta d(\theta_\#^* \alpha) + \int \varphi^\circ \circ g_\theta d(\theta_\#^* \beta) \right] d\hat{\sigma}_K(\theta).$$

In particular, this result is valid for SOT.

Proof. The proof consists in applying Lemma A.3 for (X, Y) chosen as $X = \text{supp}(\hat{\sigma}_K) \subset \mathbb{S}^{d-1}$ and

$$Y = \{\forall \theta \in \text{supp}(\hat{\sigma}_K), f_\theta : \mathbb{R} \rightarrow \mathbb{R}, g_\theta : \mathbb{R} \rightarrow \mathbb{R}, f_\theta(x) + g_\theta(y) \leq C_1(x, y)\}.$$

The functions in Y are dual potentials, and by definition are continuous for any θ . Let $\mathcal{F} : X \times Y \rightarrow \mathbb{R}$ be the functional defined as

$$\mathcal{F} : (\theta, (f_\theta)_\theta, (g_\theta)_\theta) \mapsto \int f_\theta d(\theta_\#^* \alpha) + \int g_\theta d(\theta_\#^* \beta).$$

Since the measures (α, β) have compact support, then by Lemma A.5, the supremum is attained over a subset of dual potentials of Y such that for any fixed $\theta \in X$, (f_θ, g_θ) are Lipschitz-continuous and bounded, thus uniformly equicontinuous functions (with constants independent of θ). By the Ascoli-Arzelà theorem, the set of uniformly equicontinuous functions is compact for the uniform convergence. Hence, for any $\theta \in X$, there exists a sequence of dual potentials $(f_{\theta, n}, g_{\theta, n})$ which uniformly converges to optimal dual potentials (f_θ, g_θ) (up to extraction of subsequence). Besides, we have $\text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta) = \mathcal{F}(\theta, f_\theta, g_\theta)$ and $\mathcal{F}(\theta, (f_{\theta, n})_\theta, (g_{\theta, n})_\theta) \rightarrow \text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ as $n \rightarrow +\infty$. Denote $\mathcal{F}_n(\theta) \triangleq \mathcal{F}(\theta, (f_{\theta, n})_\theta, (g_{\theta, n})_\theta)$ and $\text{OT}(\theta) \triangleq \text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$. In order to apply Lemma A.3, we need to prove that the convergence of $(\mathcal{F}_n(\theta))_{n \in \mathbb{N}^*}$ to $\text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ is uniform w.r.t. θ , i.e., $\sup_{\theta \in X} |\mathcal{F}_n(\theta) - \text{OT}(\theta)| \rightarrow 0$ as $n \rightarrow +\infty$.

First, note that for any $\theta \in X$,

$$|\mathcal{F}_n(\theta) - \text{OT}(\theta)| \leq m(\alpha) \|f_{\theta, n} - f_\theta\|_\infty + m(\beta) \|g_{\theta, n} - g_\theta\|_\infty.$$

Since for a fixed $\theta \in X$, $(f_{\theta, n}, g_{\theta, n})_{n \in \mathbb{N}^*}$ uniformly converge to (f_θ, g_θ) , this means that

$$\forall \theta \in X, \forall \varepsilon > 0, \exists N(\varepsilon, \theta), \forall n \geq N(\varepsilon, \theta), m(\alpha) \|f_{\theta, n} - f_\theta\|_\infty + m(\beta) \|g_{\theta, n} - g_\theta\|_\infty < \varepsilon.$$

Since we assume that σ is supported on a discrete set, then the cardinal of X is finite and one can define $N(\varepsilon) \triangleq \max_{\theta \in X} N(\varepsilon, \theta)$. This yields,

$$\forall \varepsilon > 0, \exists N(\varepsilon), \forall n \geq N(\varepsilon), \sup_{\theta \in X} |\mathcal{F}_n(\theta) - \text{OT}(\theta)| < \varepsilon.$$

which means that $\sup_{\theta \in X} |\mathcal{F}_n(\theta) - \text{OT}(\theta)| \rightarrow 0$, thus concludes the proof. \square

Lemma A.5. *Let $p \in [1, +\infty)$ and $C_1(x, y) = |x - y|^p$. Consider two positive measures $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d)$ whose support is such that $C_d(x, y) = \|x - y\|^p \leq R$. Then for any $\theta \in \mathbb{S}^{d-1}$, one can restrict without loss of generality the problem $\text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ as a supremum over dual potentials satisfying $f_\theta(x) + g_\theta(y) \leq C_1(x, y)$, uniformly bounded by M and uniformly L -Lipschitz, where M and L do not depend on θ .*

Proof. We adapt the proof of (Santambrogio, 2015, Proposition 1.11), and focus on showing that the uniform boundedness and Lipschitz constant are independent of $\theta \in \mathbb{S}^{d-1}$ in this setting. Here we consider the translation-invariant formulation of UOT from (Séjourné et al., 2022b), *i.e.*, $\text{UOT}(\alpha, \beta) = \sup_{f \oplus g \leq C_d} \mathcal{H}(f, g)$, where $\mathcal{H}(f, g) = \sup_{\lambda \in \mathbb{R}} \mathcal{D}(f + \lambda, g - \lambda)$. It is proved in (Séjourné et al., 2022b, Proposition 9) that the above problem has the same primal and is thus equivalent to optimize \mathcal{D} . By definition one has $\mathcal{H}(f, g) = \mathcal{H}(f + \lambda, g - \lambda)$ for any $\lambda \in \mathbb{R}$, *i.e.*, this formulation shares the same invariance as Balanced OT. Thus we can reuse all arguments from (Santambrogio, 2015, Proposition 1.11), such that for $\text{UOT}(\alpha, \beta)$, one can use the constraint $f(x) + g(y) \leq C_d(x, y)$ and the assumption $C_d(x, y) \leq R$ to prove that without loss of generality, one can restrict to potentials such that $f(x) \in [0, R]$ and $g(y) \in [-R, R]$. Furthermore if the cost satisfies in \mathbb{R}^d

$$|C_d(x, y) - C_d(x', y')| \leq L(\|x - x'\| + \|y - y'\|),$$

then one can also restrict w.l.o.g. to potentials which are L -Lipschitz. For the cost $C_d(x, y) = \|x - y\|^p$ with $p \geq 1$, this holds with constant $L = pR^{p-1}$ because the support is bounded and the gradient of C_d is radially non-decreasing.

Regarding $\text{OT}(\theta_{\#}^* \alpha, \theta_{\#}^* \beta)$, the bounds (M_θ, L_θ) could be refined by considering the dependence in $\theta \in \mathbb{S}^{d-1}$. However we prove now these constants can be upper-bounded by a finite constant independent of θ . In this setting we consider the cost

$$C_1(\theta^*(x), \theta^*(y)) = |\langle \theta, x - y \rangle|^p \leq \|\theta\|^p \|x - y\|^p \leq \|x - y\|^p,$$

by Cauchy-Schwarz inequality. Therefore, if (α, β) have supports such that $\|x - y\|^p \leq R$, then $(\theta_{\#}^* \alpha, \theta_{\#}^* \beta)$ also have supports bounded by R in \mathbb{R} . Similarly note that the derivative of $h(x) = x^p$ is non-decreasing for $p \geq 1$. Hence the cost $C_1(\theta^*(x), \theta^*(y))$ has a bounded derivative, which reads

$$p |\langle \theta, x - y \rangle|^{p-1} \leq p \|\theta\|^{p-1} \|x - y\|^{p-1} \leq p \|x - y\|^{p-1} \leq pR^{p-1}.$$

Thus on the supports of $(\theta_{\#}^* \alpha, \theta_{\#}^* \beta)$ one can also bound the Lipschitz constant of the cost $C_1(x, y) = |x - y|^p$ by the same constant L . \square

Remark: Extending Theorem 3.9. We conjecture that Theorem 3.9 also holds when σ is the uniform measures over \mathbb{S}^{d-1} , since the above holds for any $N \in \mathbb{N}^*$ and $\hat{\sigma}_N$ converges weakly* to σ . Proving this result would require that potentials (f_θ, g_θ) are also regular (*i.e.*, Lipschitz and bounded) w.r.t $\theta \in \mathbb{S}^{d-1}$. This regularity is proved in (Xi & Niles-Weed, 2022) assuming (α, β) have densities, but remains unknown for discrete measures. Since discretizing σ corresponds to the computational approach, we assume it to be discrete, so that no additional assumption than boundedness on (α, β) is required. For instance, such result remains valid for semi-discrete UOT computation.

A.3 Metric properties: Proof of Proposition 3.3 and Proposition 3.8

Proof of Proposition 3.3. Metric properties of SUOT. Symmetry and non-negativity are immediate. Assume $\text{SUOT}(\alpha, \beta) = 0$. Since σ is the uniform distribution on \mathbb{S}^{d-1} , then for any $\theta \in \mathbb{S}^{d-1}$, $\text{UOT}(\theta_{\#}^* \alpha, \theta_{\#}^* \beta) = 0$, and since UOT is assumed to be definite, then $\theta_{\#}^* \alpha = \theta_{\#}^* \beta$. By (Bogachev & Ruas, 2007, Proposition 3.8.6), this implies that α and β have the same Fourier transform. By injectivity of the Fourier transform, we conclude that $\alpha = \beta$, hence SUOT is definite. The triangle inequality results from applying

the Minkowski inequality then the triangle inequality for $\text{UOT}^{1/p}$ for $p \in [1, +\infty)$: for any $\alpha, \beta, \gamma \in \mathcal{M}_+(\mathbb{R}^d)$,

$$\begin{aligned}
& \text{SUOT}^{1/p}(\alpha, \beta) \\
&= \left(\int_{\mathbb{S}^{d-1}} \text{UOT}(\theta_{\#}^* \alpha, \theta_{\#}^* \beta) d\boldsymbol{\sigma}(\theta) \right)^{1/p} \\
&\leq \left(\int_{\mathbb{S}^{d-1}} [\text{UOT}^{1/p}(\theta_{\#}^* \alpha, \theta_{\#}^* \gamma) + \text{UOT}^{1/p}(\theta_{\#}^* \gamma, \theta_{\#}^* \beta)]^p d\boldsymbol{\sigma}(\theta) \right)^{1/p} \\
&\leq \left(\int_{\mathbb{S}^{d-1}} [\text{UOT}^{1/p}(\theta_{\#}^* \alpha, \theta_{\#}^* \gamma)]^p d\boldsymbol{\sigma}(\theta) \right)^{1/p} + \left(\int_{\mathbb{S}^{d-1}} [\text{UOT}^{1/p}(\theta_{\#}^* \gamma, \theta_{\#}^* \beta)]^p d\boldsymbol{\sigma}(\theta) \right)^{1/p} \\
&= \text{SUOT}^{1/p}(\alpha, \gamma) + \text{SUOT}^{1/p}(\gamma, \beta).
\end{aligned}$$

□

Proof of Proposition 3.8. Metric properties of USOT. Let $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d)$. Non-negativity is immediate, as USOT is defined as a program minimizing a sum of positive terms. SOT is symmetric, thus when $\varphi_1 = \varphi_2$, we obtain symmetry of the functional w.r.t. (α, β) . Assume D_φ is definite, *i.e.*, $D_\varphi(\alpha|\beta) = 0$ implies $\alpha = \beta$. Assume now that $\text{USOT}(\alpha, \beta) = 0$, and denote by (π_1, π_2) the optimal marginals attaining the infimum in (9). $\text{USOT}(\alpha, \beta) = 0$ implies that $\text{SOT}(\pi_1, \pi_2) = 0$, $D_\varphi(\pi_1|\alpha) = 0$ and $D_\varphi(\pi_2|\beta) = 0$. These three terms are definite, which yields $\alpha = \pi_1 = \pi_2 = \beta$, hence the definiteness of USOT. The Partial OT setting (*i.e.*, $D_\varphi = \rho\text{TV}$) is treated in Appendix A.7.

□

A.4 Sample complexity: Proof of Theorem 3.4

Theorem 3.4 is obtained by adapting (Nadjahi et al., 2020b, Theorems 4 and 5). We provide the detailed derivations below.

Proof of Theorem 3.4. Let $(\alpha, \beta) \in \tilde{\mathcal{M}} \times \tilde{\mathcal{M}}$ with respective empirical approximations $\hat{\alpha}_n, \hat{\beta}_n$ over n samples. By using the definition of SUOT, the triangle inequality and the assumed sample complexity of UOT for univariate measures, we show that

$$\mathbb{E} \left| \text{SUOT}(\alpha, \beta) - \text{SUOT}(\hat{\alpha}_n, \hat{\beta}_n) \right| \quad (15)$$

$$= \mathbb{E} \left| \int_{\mathbb{S}^{d-1}} \{ \text{UOT}(\theta_{\#}^* \alpha, \theta_{\#}^* \beta) - \text{UOT}(\theta_{\#}^* \hat{\alpha}_n, \theta_{\#}^* \hat{\beta}_n) \} d\boldsymbol{\sigma}(\theta) \right| \quad (16)$$

$$\leq \mathbb{E} \left\{ \int_{\mathbb{S}^{d-1}} | \text{UOT}(\theta_{\#}^* \alpha, \theta_{\#}^* \beta) - \text{UOT}(\theta_{\#}^* \hat{\alpha}_n, \theta_{\#}^* \hat{\beta}_n) | d\boldsymbol{\sigma}(\theta) \right\} \quad (17)$$

$$\leq \int_{\mathbb{S}^{d-1}} \mathbb{E} | \text{UOT}(\theta_{\#}^* \alpha, \theta_{\#}^* \beta) - \text{UOT}(\theta_{\#}^* \hat{\alpha}_n, \theta_{\#}^* \hat{\beta}_n) | d\boldsymbol{\sigma}(\theta) \quad (18)$$

$$\leq \int_{\mathbb{S}^{d-1}} \kappa(n) d\boldsymbol{\sigma}(\theta) = \kappa(n), \quad (19)$$

which completes the proof for the first setting.

Next, let $\alpha \in \tilde{\mathcal{M}}$ with corresponding empirical approximation $\hat{\alpha}_n$. Then, using the definition of SUOT, the triangle inequality (w.r.t. integral) and the assumed convergence rate in UOT,

$$\mathbb{E} | \text{SUOT}(\hat{\alpha}_n, \alpha) | \quad (20)$$

$$= \mathbb{E} \left| \int_{\mathbb{S}^{d-1}} \text{UOT}(\theta_{\#}^* \hat{\alpha}_n, \theta_{\#}^* \alpha) d\boldsymbol{\sigma}(\theta) \right| \leq \mathbb{E} \left\{ \int_{\mathbb{S}^{d-1}} | \text{UOT}(\theta_{\#}^* \hat{\alpha}_n, \theta_{\#}^* \alpha) | d\boldsymbol{\sigma}(\theta) \right\} \quad (21)$$

$$\leq \int_{\mathbb{S}^{d-1}} \mathbb{E} | \text{UOT}(\theta_{\#}^* \hat{\alpha}_n, \theta_{\#}^* \alpha) | d\boldsymbol{\sigma}(\theta) \leq \int_{\mathbb{S}^{d-1}} \xi(n) d\boldsymbol{\sigma}(\theta) = \xi(n). \quad (22)$$

□

Corollary A.6. *Assume for $\mu \in \mathcal{M}_+(\mathbb{R})$, $\mathbb{E}|\text{UOT}(\mu, \hat{\mu}_n)| \leq \xi(n)$ and that for $p \geq 1$, $\text{UOT}^{1/p}$ satisfies non-negativity, symmetry and the triangle inequality on $\mathcal{M}_+(\mathbb{R}) \times \mathcal{M}_+(\mathbb{R})$. Then, for $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^d)$,*

$$\mathbb{E} \left| \text{SUOT}^{1/p}(\alpha, \beta) - \text{SUOT}^{1/p}(\hat{\alpha}_n, \hat{\beta}_n) \right| \leq 2\xi(n)^{1/p}. \quad (23)$$

Proof. Since $\text{UOT}^{1/p}$ satisfies non-negativity, symmetry and the triangle inequality on $\mathcal{M}_+(\mathbb{R}) \times \mathcal{M}_+(\mathbb{R})$, $\text{SUOT}^{1/p}$ verifies these three metric properties on $\mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$ by Proposition 3.3, and we can derive its sample complexity as follows. For any α, β in $\mathcal{M}_+(\mathbb{R}^d)$ with respective empirical approximations $\hat{\alpha}_n, \hat{\beta}_n$, applying the triangle inequality yields for $p \in [1, +\infty)$,

$$\left| \text{UOT}^{1/p}(\alpha, \beta) - \text{UOT}^{1/p}(\hat{\alpha}_n, \hat{\beta}_n) \right| \leq \text{UOT}^{1/p}(\hat{\alpha}_n, \alpha) + \text{UOT}^{1/p}(\hat{\beta}_n, \beta). \quad (24)$$

Taking the expectation of (24) with respect to $\hat{\alpha}_n, \hat{\beta}_n$ gives,

$$\mathbb{E} \left| \text{SUOT}^{1/p}(\alpha, \beta) - \text{SUOT}^{1/p}(\hat{\alpha}_n, \hat{\beta}_n) \right| \leq \mathbb{E}|\text{SUOT}^{1/p}(\hat{\alpha}_n, \alpha)| + \mathbb{E}|\text{SUOT}^{1/p}(\hat{\beta}_n, \beta)| \quad (25)$$

$$\leq \{\mathbb{E}|\text{SUOT}(\hat{\alpha}_n, \alpha)|\}^{1/p} + \{\mathbb{E}|\text{SUOT}(\hat{\beta}_n, \beta)|\}^{1/p} \quad (26)$$

$$\leq \xi(n)^{1/p} + \xi(n)^{1/p} = 2\xi(n)^{1/p}, \quad (27)$$

where (26) is immediate if $p = 1$, and results from applying Hölder's inequality on \mathbb{S}^{d-1} if $p > 1$, and (27) follows from (22). □

A.5 Comparison of SUOT, USOT, SOT, and proof of Theorem 3.10 and Theorem 3.11

In this section, we establish several bounds to compare SUOT, USOT and SOT on the space of compactly-supported measures. We provide the detailed derivations and auxiliary lemmas needed for the proofs. The Partial OT setting (*i.e.*, $D_\varphi = \rho\text{TV}$) is treated in Appendix A.7.

Proof of Theorem 3.10. Theorem 3.10 is a direct consequence from Theorems A.7 and A.8.

Theorem A.7. *Let $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$. Then, $\text{SUOT}(\alpha, \beta) \leq \text{USOT}(\alpha, \beta)$.*

Proof. To show that $\text{SUOT}(\alpha, \beta) \leq \text{USOT}(\alpha, \beta)$, we use a sub-optimality argument. Let π be the solution $\text{USOT}(\alpha, \beta)$ and denote by (π_1, π_2) the marginals of π . For any $\theta \in \mathbb{S}^{d-1}$, denote by π_θ the solution of $\text{OT}(\theta_\#^* \pi_1, \theta_\#^* \pi_2)$. By definition of USOT, the marginals of π_θ are given by $(\theta_\#^* \pi_1, \theta_\#^* \pi_2)$. Since the sequence $(\pi_\theta)_\theta$ is suboptimal for the problem $\text{SUOT}(\alpha, \beta)$, one has

$$\text{SUOT}(\alpha, \beta) \leq \int_{\mathbb{S}^{d-1}} \left\{ \int C_1 d\pi_\theta + D_{\varphi_1}(\theta_\#^* \pi_1 | \theta_\#^* \alpha) + D_{\varphi_2}(\theta_\#^* \pi_2 | \theta_\#^* \beta) \right\} d\sigma(\theta) \quad (28)$$

$$\leq \int_{\mathbb{S}^{d-1}} \int C_1 d\pi_\theta d\sigma(\theta) + D_{\varphi_1}(\pi_1 | \alpha) + D_{\varphi_2}(\pi_2 | \beta) \quad (29)$$

$$= \text{USOT}(\alpha, \beta), \quad (30)$$

where the second inequality results from Lemma A.10, and the last equality follows from the definition of $\text{USOT}(\alpha, \beta)$. □

Theorem A.8. *Let $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d) \times \mathcal{M}_+(\mathbb{R}^d)$. Additionally, let $p \in [1, +\infty)$ and assume that for all $x, y \in \mathbb{R}^d$, $\theta \in \mathbb{S}^{d-1}$, $C_1(\theta^*(x), \theta^*(y)) \leq C_d(x, y)$. Then, $\text{USOT}(\alpha, \beta) \leq \text{UOT}(\alpha, \beta)$.*

Proof. By (Bonnotte, 2013, Proposition 5.1.3), $\text{SOT}(\mu, \nu) \leq \text{OT}(\mu, \nu)$ as $C_1(\theta^*(x), \theta^*(y)) \leq C_d(x, y)$ for all $x, y \in \mathbb{R}^d$, $\theta \in \mathbb{S}^{d-1}$. Let π be the solution of $\text{UOT}(\alpha, \beta)$ with marginals (π_1, π_2) . These marginals are

sub-optimal for $\text{USOT}(\alpha, \beta)$, we have

$$\text{USOT}(\alpha, \beta) \leq \text{SOT}(\pi_1, \pi_2) + D_{\varphi_1}(\pi_1|\alpha) + D_{\varphi_2}(\pi_2|\beta), \quad (31)$$

$$\leq \text{OT}(\pi_1, \pi_2) + D_{\varphi_1}(\pi_1|\alpha) + D_{\varphi_2}(\pi_2|\beta), \quad (32)$$

$$= \text{UOT}(\alpha, \beta), \quad (33)$$

where the last equality is obtained because π is optimal in $\text{UOT}(\alpha, \beta)$. □

Proof of Theorem 3.11.

Theorem A.9. *Let X be a compact subset of \mathbb{R}^d with radius R and consider $\alpha, \beta \in \mathcal{M}_+(\mathbb{X})$. Additionally, let $p \in [1, +\infty)$ and assume $C_1(x, y) = |x - y|^p$ for $(x, y) \in \mathbb{R}$ and $C_d(x, y) = \|x - y\|^p$ for $(x, y) \in \mathbb{R}^d$. Let $\rho > 0$ and assume $D_{\varphi_1} = D_{\varphi_2} = \rho \text{KL}$. Then, $\text{UOT}(\alpha, \beta) \leq c \text{SUOT}(\alpha, \beta)^{1/(d+1)}$, where $c = c(m(\alpha), m(\beta), \rho, R)$ is a non-decreasing function of $m(\alpha)$ and $m(\beta)$.*

Proof. We adapt the proof of (Bonnotte, 2013, Lemma 5.1.4), which establishes a bound between OT and SOT. The first step consists in bounding from above the distance between two regularized measures.

Let $\psi : \mathbb{R}^d \rightarrow \mathbb{R}_+$ be a smooth and radial function such that $\text{supp}(\psi) \subseteq B_d(\mathbf{0}, 1)$ and $\int_{\mathbb{R}^d} \psi(x) d\text{Leb}(x) = 1$. Let $\psi_\lambda(x) = \lambda^{-d} \psi(x/\lambda)$. For any function f defined on \mathbb{R}^s ($s \geq 1$), denote by $\mathcal{F}[f]$ the Fourier transform of f defined for $x \in \mathbb{R}^s$ as $\mathcal{F}[f](x) = \int_{\mathbb{R}^s} f(w) e^{-i\langle w, x \rangle} dw$. Let $\alpha_\lambda = \alpha * \varphi_\lambda$ and $\beta_\lambda = \beta * \varphi_\lambda$ where $*$ is the convolution operator. Let (f, g) such that $f \oplus g \leq C_d$. By using the isometry properties of the Fourier transform and the definition of ψ_λ , then representing the variables with polar coordinates, we have

$$\int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha_\lambda(x) = \int_{\mathbb{R}^d} \mathcal{F}[\varphi^\circ \circ f](w) \mathcal{F}[\alpha](w) \mathcal{F}[\psi](\lambda w) dw \quad (34)$$

$$= \int_{\mathbb{S}^{d-1}} \int_0^{+\infty} \mathcal{F}[\varphi^\circ \circ f](r\theta) \mathcal{F}[\alpha](r\theta) \mathcal{F}[\psi](\lambda r) r^{d-1} dr d\theta. \quad (35)$$

Since $\varphi^\circ \circ f$ is a real-valued function, $\mathcal{F}[\varphi^\circ \circ f]$ is an even function, then

$$\int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha_\lambda(x) \quad (36)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \mathcal{F}[\varphi^\circ \circ f](r\theta) \mathcal{F}[\alpha](r\theta) \mathcal{F}[\psi](\lambda r) |r|^{d-1} dr d\theta \quad (37)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \mathcal{F}[\varphi^\circ \circ f](r\theta) \mathcal{F}[\theta_\#^* \alpha](r) \mathcal{F}[\psi](\lambda r) |r|^{d-1} dr d\theta \quad (38)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \mathcal{F}[\varphi^\circ \circ f](r\theta) \left(\int_{-R}^R e^{-iru} d\theta_\#^* \alpha(u) \right) \mathcal{F}[\psi](\lambda r) |r|^{d-1} dr d\theta \quad (39)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \left(\int_{\mathbb{R}^d} \int_{-R}^R \varphi^\circ(f(x)) e^{-ir\langle u, \theta, x \rangle} d\theta_\#^* \alpha(u) \right) \mathcal{F}[\psi](\lambda r) |r|^{d-1} dx dr d\theta. \quad (40)$$

Equation 38 follows from the property of push-forward measures, (39) results from the definition of the Fourier transform and $u \in [-R, R]$, and (40) results from the definition of the Fourier transform and Fubini's theorem. By making a change of variables (x becomes $x - u\theta$), we obtain

$$\int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha_\lambda(x) \quad (41)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \int_{\mathbb{R}^d} \int_{-R}^R \varphi^\circ(f(x - u\theta)) e^{-ir\langle \theta, x \rangle} d\theta_\#^* \alpha(u) \mathcal{F}[\psi](\lambda r) |r|^{d-1} dx dr d\theta \quad (42)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \int_{B_d(\mathbf{0}, 2R+\lambda)} \int_{-R}^R \varphi^\circ(f(x - u\theta)) e^{-ir\langle \theta, x \rangle} d\theta_\#^* \alpha(u) \mathcal{F}[\psi](\lambda r) |r|^{d-1} dx dr d\theta, \quad (43)$$

where (43) follows from $\text{supp}(\alpha) \subseteq B_d(\mathbf{0}, R)$. Indeed, this implies that $\text{supp}(\alpha_\lambda) \subseteq B_d(\mathbf{0}, R + \lambda)$, thus the domain of $x \mapsto \varphi^\circ \circ f(x - u\theta)$ is contained in $B_d(\mathbf{0}, 2R + \lambda)$.

Similarly, one can show that

$$\int_{\mathbb{R}^d} \varphi^\circ(g(y)) d\beta_\lambda(y) \quad (44)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \int_{B_d(\mathbf{0}, 2R + \lambda)} \int_{-R}^R \varphi^\circ(g(y - u\theta)) e^{-ir\langle \theta, y \rangle} d\theta_\#^* \beta(u) \mathcal{F}[\psi](\lambda r) |r|^{d-1} dy dr d\theta. \quad (45)$$

By (43) and (45), we obtain

$$\int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha_\lambda(x) + \int_{\mathbb{R}^d} \varphi^\circ(g(y)) d\beta_\lambda(y) \quad (46)$$

$$= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{R}} \int_{B_d(\mathbf{0}, 2R + \lambda)} \left\{ \int_{-R}^R \varphi^\circ(f(x - u\theta)) d\theta_\#^* \alpha(u) + \int_{-R}^R \varphi^\circ(g(x - u\theta)) d\theta_\#^* \beta(u) \right\} e^{-ir\langle \theta, x \rangle} \mathcal{F}[\psi](\lambda r) |r|^{d-1} dx dr d\theta, \quad (47)$$

and,

$$\left| \int_{B_d(\mathbf{0}, 2R + \lambda)} \left\{ \int_{-R}^R \varphi^\circ(f(x - u\theta)) d\theta_\#^* \alpha(u) + \int_{-R}^R \varphi^\circ(g(x - u\theta)) d\theta_\#^* \beta(u) \right\} e^{-ir\langle \theta, x \rangle} dx \right| \quad (48)$$

$$\leq \int_{B_d(\mathbf{0}, 2R + \lambda)} \text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta) |e^{-ir\langle \theta, x \rangle}| dx \quad (49)$$

$$\leq (2R + \lambda)^d \text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta), \quad (50)$$

where (49) is obtained by taking the supremum of (48) over the set of potentials (\tilde{f}, \tilde{g}) such that for $u \in [-R, R]$, $\exists(x, \theta) \in B_d(\mathbf{0}, 2R + \lambda) \times \mathbb{S}^{d-1}$, $\tilde{f}(u) = f(x - u\theta)$, $\tilde{g}(u) = g(x - u\theta)$, which is included in the set of potentials (f', g') s.t. $f' : \mathbb{R} \rightarrow \mathbb{R}$, $g' : \mathbb{R} \rightarrow \mathbb{R}$ and $f' \oplus g' \leq C_1$. Therefore,

$$\int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha_\lambda(x) + \int_{\mathbb{R}^d} \varphi^\circ(g(y)) d\beta_\lambda(y) \leq \frac{1}{2} (2R + \lambda)^d \mathcal{A}(\mathbb{S}^{d-1}) \int_{\mathbb{S}^{d-1}} \text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta) d\sigma(\theta) \int_{\mathbb{R}} \lambda^{-d} |\mathcal{F}[\psi](r)| |r|^{d-1} |dr| \quad (51)$$

$$\leq c(2R + \lambda)^d \lambda^{-d} \text{SUOT}(\alpha, \beta), \quad (52)$$

where $c = \frac{1}{2} \mathcal{A}(\mathbb{S}^{d-1}) \int_{\mathbb{R}} |\mathcal{F}[\psi](r)| |r|^{d-1} |dr|$. We deduce from the dual formulation of UOT (3) and (49) that,

$$\text{UOT}(\alpha_\lambda, \beta_\lambda) \leq c(2R + \lambda)^d \lambda^{-d} \text{SUOT}(\alpha, \beta). \quad (53)$$

The last step of the proof consists in relating $\text{UOT}(\alpha_\lambda, \beta_\lambda)$ with $\text{UOT}(\alpha, \beta)$. For any (f, g) such that $f \oplus g \leq C_d$, we have

$$\int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha(x) + \int_{\mathbb{R}^d} \varphi^\circ(g(y)) d\beta(y) - \text{UOT}(\alpha_\lambda, \beta_\lambda) \quad (54)$$

$$\leq \int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha(x) + \int_{\mathbb{R}^d} \varphi^\circ(g(x)) d\beta(x) - \int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha_\lambda(x) - \int_{\mathbb{R}^d} \varphi^\circ(g(y)) d\beta_\lambda(y) \quad (55)$$

$$\leq \int_{\mathbb{R}^d} \{\varphi^\circ(f(x)) - \psi_\lambda * \varphi^\circ(f(x))\} d\alpha(x) + \int_{\mathbb{R}^d} \{\varphi^\circ(g(y)) - \psi_\lambda * \varphi^\circ(g(y))\} d\beta(y). \quad (56)$$

For $x \in \mathbb{R}^d$,

$$\varphi^\circ(f(x)) - \psi_\lambda * \varphi^\circ(f(x)) = \lambda^{-d} \int_{\mathbb{R}^d} (\varphi^\circ(f(x)) - \varphi^\circ(f(y))) \psi\left(\frac{x-y}{\lambda}\right) dy \quad (57)$$

$$\leq \lambda^{-d} \int_{\mathbb{R}^d} |\varphi^\circ(f(x)) - \varphi^\circ(f(y))| \psi\left(\frac{x-y}{\lambda}\right) dy, \quad (58)$$

Since $D_\varphi = \rho\text{KL}$, then for $z \in \mathbb{R}$, $\varphi^\circ(z) = \rho(1 - e^{-z/\rho})$, so for $(x, y) \in \mathbb{R}^d \times \mathbb{R}^d$,

$$\varphi^\circ(f(x)) - \varphi^\circ(f(y)) = \rho(e^{-f(y)/\rho} - e^{-f(x)/\rho}) \quad (59)$$

By Lemma A.13, the potentials (f, g) are bounded by constants depending on $m(\alpha), m(\beta)$. Therefore, we can bound (59) as follows.

$$|\varphi^\circ(f(x)) - \varphi^\circ(f(y))| \leq e^{-\lambda^*/\rho} \|x - y\|. \quad (60)$$

We thus derive the following upper-bound on (58).

$$\varphi^\circ(f(x)) - \psi_\lambda * \varphi^\circ(f(x)) \leq \lambda^{-d} e^{-\lambda^*/\rho} \int_{\mathbb{R}^d} \|x - y\| \psi\left(\frac{x-y}{\lambda}\right) dy \quad (61)$$

$$\leq \lambda^{-d+1} e^{-\lambda^*/\rho} \int_{\mathbb{R}^d} \frac{\|x - y\|}{\lambda} \psi\left(\frac{x-y}{\lambda}\right) dy \quad (62)$$

By doing the change of variables $z = (y - x)/\lambda$ and using the fact that ψ is a radial function, we obtain

$$\varphi^\circ(f(x)) - \psi_\lambda * \varphi^\circ(f(x)) \leq \lambda e^{-\lambda^*/\rho} \int_{\mathbb{R}^d} \|z\| \psi(z) dz. \quad (63)$$

Similarly, using the bounds on g in Lemma A.13, one can show that

$$|\varphi^\circ(g(x)) - \varphi^\circ(g(y))| \leq e^{(\lambda^*+R)/\rho} \|x - y\|, \quad (64)$$

therefore,

$$\varphi^\circ(g(x)) - \psi_\lambda * \varphi^\circ(g(x)) \leq \lambda e^{(\lambda^*+R)/\rho} \int_{\mathbb{R}^d} \|z\| \psi(z) dz. \quad (65)$$

We conclude that,

$$\int_{\mathbb{R}^d} \varphi^\circ(f(x)) d\alpha(x) + \int_{\mathbb{R}^d} \varphi^\circ(g(y)) d\beta(y) - \text{UOT}(\alpha_\lambda, \beta_\lambda) \leq \left(e^{-\lambda^*/\rho} + e^{(\lambda^*+R)/\rho} \right) \lambda M_1(\psi), \quad (66)$$

where $M_1(\psi) \triangleq \int_{\mathbb{R}^d} \|z\| \psi(z) dz$. Taking the supremum on both sides over (f, g) such that $f \oplus g \leq C_d$ yields,

$$\text{UOT}(\alpha, \beta) - \text{UOT}(\alpha_\lambda, \beta_\lambda) \leq \left(e^{-\lambda^*/\rho} + e^{(\lambda^*+R)/\rho} \right) \lambda M_1(\psi). \quad (67)$$

Finally, by combining (53) with the above inequality, we obtain

$$\text{UOT}(\alpha, \beta) \leq \left(e^{-\lambda^*/\rho} + e^{(\lambda^*+R)/\rho} \right) \lambda M_1(\psi) + c(2R + \lambda)^d \lambda^{-d} \text{SUOT}(\alpha, \beta) \quad (68)$$

$$\leq c' \lambda (1 + (2R + \lambda)^d \lambda^{-(d+1)}) \text{SUOT}(\alpha, \beta), \quad (69)$$

where c' is a constant satisfying $c' \geq c$ and $c' \geq (e^{-\lambda^*/\rho} + e^{(\lambda^*+R)/\rho}) M_1(\psi)$. By choosing $\lambda = R^{d/(d+1)} \text{SUOT}(\alpha, \beta)^{1/(d+1)}$, (69) becomes

$$\text{UOT}(\alpha, \beta) \leq c' R^{d/(d+1)} \text{SUOT}(\alpha, \beta)^{1/(d+1)} (1 + (2R + \lambda)^d R^{-d}). \quad (70)$$

We conclude using that $\text{SUOT}(\alpha, \beta)$ is bounded from above. Indeed, $\text{SUOT}(\alpha, \beta) \leq \rho(m(\alpha) + m(\beta))$ since on the one hand, π is suboptimal in (3) thus $\text{UOT}(\alpha, \beta) \leq \rho(m(\alpha) + m(\beta))$, and on the other hand, $m(\alpha) = m(\theta_\#^* \alpha)$ for any $\theta \in \mathbb{S}^{d-1}$. This yields $\lambda \leq R^{d/(d+1)} \rho^{1/(d+1)} (m(\alpha) + m(\beta))^{1/(d+1)}$, hence

$$\text{UOT}(\alpha, \beta) \leq c(m(\alpha), m(\beta), \rho, R) \text{SUOT}(\alpha, \beta)^{1/(d+1)}, \quad (71)$$

where $c(m(\alpha), m(\beta), \rho, R)$ is a non-decreasing function of $m(\alpha)$ and $m(\beta)$.

□

Additional Lemmas.

Lemma A.10. For any $\theta \in \mathbb{S}^{d-1}$ and $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^d)$, $D_\varphi(\theta_\#^* \alpha | \theta_\#^* \beta) \leq D_\varphi(\alpha | \beta)$.

Proof. For $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^s)$ with $s \geq 1$, the dual characterization of φ -divergences reads (Liero et al., 2018, Theorem 2.7)

$$D_\varphi(\alpha | \beta) = \sup_{f \in \mathcal{E}(\mathbb{R}^s)} \int_{\mathbb{R}^s} \varphi^\circ(f(x)) d\beta(x) - \int_{\mathbb{R}^s} f(x) d\alpha(x),$$

where $\mathcal{E}(\mathbb{R}^s)$ denotes the space of lower semi-continuous functions from \mathbb{R}^s to $\mathbb{R} \cup \{+\infty\}$. Therefore, for any $\theta \in \mathbb{S}^{d-1}$ and $\alpha, \beta \in \mathcal{M}_+(\mathbb{R}^d)$,

$$D_\varphi(\theta_\#^* \alpha | \theta_\#^* \beta) = \sup_{f \in \mathcal{E}(\mathbb{R})} \int_{\mathbb{R}} \varphi^\circ(f(t)) d(\theta_\#^* \beta)(t) - \int_{\mathbb{R}} f(t) d(\theta_\#^* \alpha)(t) \quad (72)$$

$$= \sup_{g: \mathbb{R}^d \rightarrow \mathbb{R} \text{ s.t. } \exists f \in \mathcal{E}(\mathbb{R}), g = f \circ \theta^*} \int_{\mathbb{R}^d} \varphi^\circ(g(x)) d\beta(x) - \int_{\mathbb{R}^d} g(x) d\alpha(x) \quad (73)$$

where (73) results from the definition of push-forward measures. We conclude the proof by observing that the supremum in (73) is taken over a subset of $\mathcal{E}(\mathbb{R}^d)$. □

Lemma A.11. (Santambrogio, 2015, Proposition 1.11) Let $p \in [1, +\infty)$ and assume $C_d(x, y) = \|x - y\|^p$. Let α, β with compact support, such that $C_d(x, y) \leq R^p$ for $(x, y) \in \text{supp}(\alpha) \times \text{supp}(\beta)$. Then without loss of generality the dual potentials (f, g) of $\text{UOT}(\alpha, \beta)$ satisfy $f(x) \in [0, R]$ and $g(y) \in [-R, R]$.

Lemma A.12. (Séjourné et al., 2022b, Proposition 2) Define the translation-invariant dual formulation

$$\text{UOT}(\alpha, \beta) = \sup_{f \oplus g \leq C_d} \sup_{\lambda \in \mathbb{R}} \int \varphi_1^\circ(f + \lambda) d\alpha + \int \varphi_2^\circ(g - \lambda) d\beta. \quad (74)$$

Let $\rho > 0$ and assume $D_{\varphi_1} = D_{\varphi_2} = \rho \text{KL}$. Take optimal potentials (f, g) in (74). Then optimal potentials in (3) are given by $(f + \lambda^*(f, g), g - \lambda^*(f, g))$, where the optimal translation λ^* reads

$$\lambda^*(f, g) \triangleq \frac{1}{2} \left[S_\rho^\beta(g) - S_\rho^\alpha(f) \right], \quad S_\rho^\alpha(f) \triangleq -\rho \log \int e^{-f/\rho} d\alpha,$$

and we call $S_\rho^\alpha(f)$ the soft-minimum of f . When $m(\alpha) = 1$ and $m \leq f(x) \leq M$, then $m \leq S_\rho^\alpha(f) \leq M$.

Lemma A.13. Assume (α, β) have compact support such that, for $(x, y) \in \text{supp}(\alpha) \times \text{supp}(\beta)$, $C(x, y) \leq R$. Then, without loss of generality, one can restrict the optimization of the dual formulation (3) of $\text{UOT}(\alpha, \beta)$ over the set of potentials satisfying for $(x, y) \in \text{supp}(\alpha) \times \text{supp}(\beta)$,

$$f(x) \in [\lambda^*, \lambda^* + R], \quad g(y) \in [-\lambda^* - R, -\lambda^* + R],$$

where $\lambda^* \in [-R + \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}, \frac{R}{2} + \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}]$. In particular, one has

$$f(x) \in [-R + \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}, \frac{3R}{2} + \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}], \quad g(y) \in [-\frac{3R}{2} - \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}, 2R - \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}]$$

Proof. Consider the translation-invariant dual formulation (74): if (f, g) are optimal, then for any $\lambda \in \mathbb{R}$, $(f + \lambda, g - \lambda)$ are also optimal. We leverage the structure of the dual constraint $f \oplus g \leq C_d$ with Lemma A.11. Since for $(x, y) \in \text{supp}(\alpha) \times \text{supp}(\beta)$, $C_d(x, y) \leq R$, then without loss of generality, $f(x) \in [0, R]$ and $g(y) \in [-R, R]$. The potentials (f, g) are optimal for the translation-invariant dual energy, and we need a bound for the original dual functional (3). To this end, we leverage Lemma A.12 to compute the optimal

translation, such that $(f, g) = (f + \lambda^*(f, g), g - \lambda^*(f, g))$. Let $\bar{\alpha} = \alpha/m(\alpha)$ and $\bar{\beta} = \beta/m(\beta)$ be the normalized probability measures. The translation can be written as,

$$\lambda^*(f, g) = \frac{1}{2} \left[S_{\rho}^{\bar{\beta}}(g) - S_{\rho}^{\bar{\alpha}}(f) \right] + \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}, \quad (75)$$

where the functional S_{ρ}^{α} is defined in Lemma A.12. Since $\bar{\alpha}$ and $\bar{\beta}$ are probability measures, then by (Genevay et al., 2019, Proposition 1), $f(x) \in [0, R]$ and $g(x) \in [-R, R]$ respectively imply $S_{\rho}^{\bar{\alpha}}(f) \in [0, R]$ and $S_{\rho}^{\bar{\beta}}(g) \in [-R, R]$. Combining these bounds on $S_{\rho}^{\bar{\alpha}}(f)$, $S_{\rho}^{\bar{\beta}}(g)$ with the expression of $\lambda^*(f, g)$, (75) yields the desired bounds on the optimal potentials (f, g) of the dual formulation (3). \square

A.6 Metrizing weak* convergence: Proof of Theorem 3.12

The Kullback-Leibler setting is treated here. The Partial OT setting (*i.e.*, $D_{\varphi} = \rho\text{TV}$) is treated in Appendix A.7.

Proof. Let (α_n) be a sequence of measures in $\mathcal{M}_+(\mathbb{X})$ and $\alpha \in \mathcal{M}_+(\mathbb{X})$, where $\mathbb{X} \subset \mathbb{R}^d$ is compact with radius $R > 0$. First, we assume that $\alpha_n \rightharpoonup \alpha$. Then, by (Liero et al., 2018, Theorem 2.25), under our assumptions, $\alpha_n \rightharpoonup \alpha$ is equivalent to $\lim_{n \rightarrow +\infty} \text{UOT}(\alpha_n, \alpha) = 0$. This implies that $\lim_{n \rightarrow +\infty} \text{SUOT}(\alpha_n, \alpha) = 0$ and $\lim_{n \rightarrow +\infty} \text{USOT}(\alpha_n, \alpha) = 0$, since by Theorem 3.11 and non-negativity of SUOT (Proposition 3.3),

$$0 \leq \text{SUOT}(\alpha_n, \alpha) \leq \text{USOT}(\alpha_n, \alpha) \leq \text{UOT}(\alpha_n, \alpha).$$

Conversely, assume either that $\lim_{n \rightarrow +\infty} \text{SUOT}(\alpha_n, \alpha) = 0$ or $\lim_{n \rightarrow +\infty} \text{USOT}(\alpha_n, \alpha) = 0$. First assume there exists $M > 0$ such that for large enough $n \in \mathbb{N}^*$, $m(\alpha_n) \leq M$, then by Theorem 3.11, there exists $c > 0$ such that $\text{UOT}(\alpha_n, \alpha) \leq c(\text{SUOT}(\alpha_n, \alpha))^{1/(d+1)}$. Since c does not depend on the masses $(m(\alpha_n), m(\alpha))$, it does not depend on n . By Theorem 3.11, it yields metric equivalence between SUOT, USOT and UOT, thus $\lim_{n \rightarrow +\infty} \text{UOT}(\alpha_n, \alpha) = 0$. By (Liero et al., 2018, Theorem 2.25), we eventually obtain $\alpha_n \rightharpoonup \alpha$, which is the desired result.

The remaining step thus consists in proving that the sequence of masses $(m(\alpha_n))_{n \in \mathbb{N}^*}$ is indeed uniformly bounded by $M > 0$ for large enough n . Note that for any $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d)$, one has $\text{UOT}(\alpha, \beta) \geq \rho(\sqrt{m(\alpha)} - \sqrt{m(\beta)})^2$. Indeed one has $\text{UOT}(\alpha, \beta) \geq \mathcal{D}(\lambda, -\lambda)$, where \mathcal{D} denotes the dual functional (3) and $\lambda = \frac{\rho}{2} \log \frac{m(\alpha)}{m(\beta)}$. Note that the pair $(\lambda, -\lambda)$ are feasible dual potentials for the constraint $f \oplus g \leq C_d$, because the cost C_d is positive in our setting. The property of push-forwards measures means that for any $\theta \in \mathbb{S}^{d-1}$, one has $m(\theta_{\sharp}^* \alpha) = m(\alpha)$. Therefore, we obtain the following bounds for n large enough.

$$\begin{aligned} \text{USOT}(\alpha_n, \alpha) &\geq \text{SUOT}(\alpha_n, \alpha) \geq \int_{\mathbb{S}^{d-1}} \rho \left(\sqrt{m(\theta_{\sharp}^* \alpha_n)} - \sqrt{m(\theta_{\sharp}^* \alpha)} \right)^2 d\sigma(\theta), \\ &= \rho(\sqrt{m(\alpha_n)} - \sqrt{m(\alpha)})^2. \end{aligned}$$

Hence, $\lim_{n \rightarrow +\infty} \text{SUOT}(\alpha_n, \alpha) = 0$ or $\lim_{n \rightarrow +\infty} \text{USOT}(\alpha_n, \alpha) = 0$ implies $\lim_{n \rightarrow +\infty} m(\alpha_n) = m(\alpha)$. In other terms the mass of sequence converges and is thus uniformly bounded for large enough n . Since we proved that $m(\alpha_n) < M$ and $m(\alpha)$ is finite, it ends the proof. \square

A.7 Properties of sliced partial OT

We provide in this subsection the proofs of Proposition 3.3, Theorems 3.11 and 3.12 for the setting of sliced partial OT. To this end, we rely on a formulation for SUOT and USOT when $D_{\varphi_1} = D_{\varphi_2} = \rho\text{TV}$, which we prove below. Equation 76 is proved in (Piccoli & Rossi, 2014) and can then be applied to SUOT: we include it for completeness. Equation 77 is our contribution and is specific to USOT.

Lemma A.14. Let $\rho > 0$ and assume $D_{\varphi_1} = D_{\varphi_2} = \rho\text{TV}$ and $C_d(x, y) = \|x - y\|$. Then, for any $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{R}^d)$,

$$\text{UOT}(\alpha, \beta) = \sup_{f \in \mathcal{E}} \int f(x) d(\alpha - \beta)(x), \quad (76)$$

where $\mathcal{E} = \{f : \mathbb{R}^d \rightarrow \mathbb{R}, \|f\|_{Lip} \leq 1, \|f\|_\infty \leq \rho\}$, $\|f\|_\infty \triangleq \sup_{x \in \mathbb{R}^d} |f(x)|$ and $\|f\|_{Lip} \triangleq \sup_{(x,y) \in \mathbb{R}^d} \frac{|f(x) - f(y)|}{C_d(x,y)}$.

Furthermore, for $C_1(x, y) = |x - y|$ and an empirical approximation $\hat{\sigma}_N = \frac{1}{N} \sum_{i=1}^N \delta_{\theta_i}$ of σ , one has

$$\text{USOT}(\alpha, \beta) = \sup_{(f_\theta) \in \mathcal{E}} \int_{\mathbb{R}^d} \left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\hat{\sigma}_N(\theta) \right) d(\alpha - \beta)(x), \quad (77)$$

where

$$\mathcal{E} = \{\forall \theta \in \text{supp}(\hat{\sigma}_N), f_\theta : \mathbb{R} \rightarrow \mathbb{R}, \|f_\theta\|_{Lip} \leq 1, \|\int_{\mathbb{S}^{d-1}} f_\theta \circ \theta^* d\hat{\sigma}_N(\theta)\|_\infty \leq \rho\},$$

and the Lipschitz norm here is defined w.r.t. C_1 as $\|f\|_{Lip} \triangleq \sup_{(x,y) \in \mathbb{R}^d} \frac{|f(x) - f(y)|}{C_1(x,y)}$

Proof. We start with the formulation of Equation 3 and Theorem 3.9. For USOT one has

$$\begin{aligned} \text{USOT}(\alpha, \beta) = & \sup_{f_\theta(\cdot) \oplus g_\theta(\cdot) \leq C_1} \int \varphi_1^\circ \left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\sigma_N(\theta) \right) d\alpha(x) \\ & + \int \varphi_2^\circ \left(\int_{\mathbb{S}^{d-1}} g_\theta(\theta^*(y)) d\sigma_N(\theta) \right) d\beta(y). \end{aligned}$$

When $D_\varphi = \rho\text{TV}$, the function φ° reads $\varphi^\circ(x) = x$ for $x \in [-\rho, \rho]$, $\varphi^\circ(x) = \rho$ when $x \geq \rho$, and $\varphi^\circ(x) = -\infty$ otherwise. Noting $f_{avg}(x) = \int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\sigma_N(\theta)$ and $g_{avg}(x) = \int_{\mathbb{S}^{d-1}} g_\theta(\theta^*(x)) d\sigma_N(\theta)$. This formula on φ° imposes $f_{avg}(x) \geq -\rho$ and $g_{avg}(x) \geq -\rho$. Furthermore, since we perform a supremum w.r.t. (f_{avg}, g_{avg}) where φ° attains a plateau, then without loss of generality, we can impose the constraint $f_{avg}(x) \leq \rho$ and $g_{avg}(x) \geq \rho$, as it will have no impact on the optimal dual functional value. Thus we have that $\|f_{avg}\|_\infty \leq \rho$ and $\|g_{avg}\|_\infty \leq \rho$. To obtain the Lipschitz property, we use the constraint that $f_\theta(\cdot) \oplus g_\theta(\cdot) \leq C_1$ for any $\theta \in \text{supp}(\sigma_N)$, as well as (Santambrogio, 2015, Proposition 3.1). Thus by using c-transform for the cost $C_1(x, y) = |x - y|$, we can take w.l.o.g $f_\theta(\cdot) = -g_\theta(\cdot)$ with $f_\theta(\cdot)$ a 1-Lipschitz function. Thus w.l.o.g we can perform the supremum over $(f_\theta)_\theta \in \mathcal{E}$, and rephrase the functional as desired, since we have that $\varphi^\circ(f_{avg}) = f_{avg}$.

The proof for UOT is exactly the same, except that our inputs are (f, g) instead of (f_θ, g_θ) . □

We can now prove Proposition 3.3, Theorems 3.11 and 3.12 in the setting of sliced Partial OT. All those results are summarized in the following statement.

Theorem A.15. (Properties of Sliced Partial OT) Assume $C_1(x, y) = |x - y|$ and $D_{\varphi_1} = D_{\varphi_2} = \rho\text{TV}$. Then, USOT satisfies the triangle inequality. Additionally, for any $(\alpha, \beta) \in \mathcal{M}_+(\mathbb{X})$ where $\mathbb{X} \subset \mathbb{R}^d$ is compact with radius R , $\text{UOT}(\alpha, \beta) \leq c(\rho, R) \text{SUOT}(\alpha, \beta)^{1/(d+1)}$, and USOT and SUOT both metrize the weak* convergence.

Proof of Sliced Partial OT properties. First we prove that in that setting USOT is a metric. Reusing Lemma A.14, we have that for any measures (α, β, γ)

$$\begin{aligned} \text{USOT}(\alpha, \gamma) &= \sup_{(f_\theta)_{\theta \in \mathcal{E}}} \int \left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\sigma_N \right) d(\alpha - \gamma)(x) \\ &= \sup_{(f_\theta)_{\theta \in \mathcal{E}}} \int \left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\sigma_N \right) d(\alpha - \beta + \beta - \gamma)(x) \\ &\leq \sup_{(f_\theta)_{\theta \in \mathcal{E}}} \int \left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\sigma_N \right) d(\alpha - \beta)(x) \\ &\quad + \sup_{(f_\theta)_{\theta \in \mathcal{E}}} \int \left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\sigma_N \right) d(\beta - \gamma)(x) \\ &= \text{USOT}(\alpha, \beta) + \text{USOT}(\beta, \gamma). \end{aligned}$$

Note that reusing Lemma A.14, we have that SUOT is a sliced integral probability metric over the space of bounded and Lipschitz functions. More precisely, we satisfy the assumptions of (Nadjahi et al., 2020b, Theorem 3), so that one has $\text{UOT}(\alpha, \beta) \leq c(\rho, R)(\text{SUOT}(\alpha, \beta))^{1/(d+1)}$.

To prove that USOT and SUOT metrize the weak* convergence, the proof is very similar to that of Theorem 3.12 detailed above. Assuming that $\alpha_n \rightarrow \alpha$ implies $\text{SUOT}(\alpha_n, \alpha) \rightarrow 0$ and $\text{USOT}(\alpha_n, \alpha) \rightarrow 0$ is already proved in Appendix A.6. To prove the converse, the proof is also the same, *i.e.*, we use the property that SUOT, USOT and UOT are equivalent metrics, which holds as we assumed that supports of (α, β) are compact in a ball of radius R . Note that since the bound $\text{UOT}(\alpha, \beta) \leq c(\rho, R)(\text{SUOT}(\alpha, \beta))^{1/(d+1)}$ holds independently of the measure's masses, we do not need to uniformly bound $m(\alpha_n)$, compared to the KL setting of Theorem 3.12.

□

B Additional details for Section 4

B.1 Postponed Proofs for Section 4

Proof of Proposition 4.1. Our goal is to compute the first order variation of the SUOT functional. Given that $\text{SUOT}(\alpha, \beta) = \int_{\mathbb{S}^{d-1}} \text{UOT}(\theta_\#^* \alpha, \theta_\#^* \beta) d\sigma(\theta)$, one can apply Proposition B.1 slice-wise. Since measures are assumed to have compact support, one can apply the dominated convergence theorem and differentiate under the integral sign. Furthermore, the translation-invariant formulation in the setting of SUOT reads

$$\text{SUOT}(\alpha, \beta) = \int_{\mathbb{S}^{d-1}} \sup_{f_\theta \oplus g_\theta \leq C_1} \left[\sup_{\lambda_\theta \in \mathbb{R}} \int \varphi^\circ(f_\theta(\cdot) + \lambda_\theta) d\theta_\#^* \alpha \right. \quad (78)$$

$$\left. + \int \varphi^\circ(g_\theta(\cdot) - \lambda_\theta) d\theta_\#^* \beta \right], \quad (79)$$

In the setting where φ° is smooth and strictly concave (such as $D_\varphi = \rho\text{KL}$), there always exists a unique optimal λ_θ^* . Furthermore, one can apply the envelope theorem such that the Fréchet differential w.r.t. to a perturbation (r_θ, s_θ) of (f_θ, g_θ) reads

$$\int_{\mathbb{S}^{d-1}} \left[\int r_\theta(\cdot) \times \nabla \varphi^\circ(f_\theta(\cdot) + \lambda_\theta^*(f_\theta, g_\theta)) d\theta_\#^* \alpha \right. \quad (80)$$

$$\left. + \int s_\theta(\cdot) \times \nabla \varphi^\circ(g_\theta(\cdot) - \lambda_\theta^*(f_\theta, g_\theta)) d\theta_\#^* \beta \right] \quad (81)$$

Setting

$$\alpha_\theta = \nabla \varphi^\circ(f_\theta(\cdot) + \lambda_\theta^*(f_\theta, g_\theta)) \alpha, \quad \beta_\theta = \nabla \varphi^\circ(g_\theta(\cdot) - \lambda_\theta^*(f_\theta, g_\theta)) \beta,$$

yields the desired result, *i.e.*, the first order variation is

$$\int_{\mathbb{S}^{d-1}} \left[\int r_\theta(\cdot) d(\theta_\#^* \alpha_\theta) + \int s_\theta(\cdot) d(\theta_\#^* \beta_\theta) \right]. \quad (82)$$

□

Proof of Proposition 4.2. Our goal is to compute the first order variation of the USOT functional. First, we leverage Theorem 3.9 such that USOT reads

$$\text{USOT}(\alpha, \beta) = \sup_{f_\theta(\cdot) \oplus g_\theta(\cdot) \leq C_1} \int \varphi_1^\circ \left(\int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\hat{\sigma}_K(\theta) \right) d\alpha(x) \quad (83)$$

$$+ \int \varphi_2^\circ \left(\int_{\mathbb{S}^{d-1}} g_\theta(\theta^*(y)) d\hat{\sigma}_K(\theta) \right) d\beta(y) \quad (84)$$

$$= \sup_{f_\theta(\cdot) \oplus g_\theta(\cdot) \leq C_1} \int \varphi_1^\circ(f_{avg}(x)) d\alpha(x) + \int \varphi_2^\circ(g_{avg}(y)) d\beta(y), \quad (85)$$

where

$$f_{avg}(x) = \int_{\mathbb{S}^{d-1}} f_\theta(\theta^*(x)) d\hat{\sigma}_K(\theta), \quad g_{avg}(y) = \int_{\mathbb{S}^{d-1}} g_\theta(\theta^*(y)) d\hat{\sigma}_K(\theta).$$

From this, we derive the translation-invariant formulation as follows.

$$\text{USOT}(\alpha, \beta) = \sup_{f_\theta(\cdot) \oplus g_\theta(\cdot) \leq C_1} \sup_{\lambda \in \mathbb{R}} \int \varphi_1^\circ(f_{avg}(x) + \lambda) d\alpha(x) \quad (86)$$

$$+ \int \varphi_2^\circ(g_{avg}(y) - \lambda) d\beta(y), \quad (87)$$

For smooth and strictly concave φ° , there exists a unique $\lambda^*(f_{avg}, g_{avg})$ attaining the supremum. Furthermore, one can apply the envelope theorem and differentiate under the integral sign (since the support is compact). Consider perturbations $(r_\theta(\cdot), s_\theta(\cdot))$ of $(f_\theta(\cdot), g_\theta(\cdot))$. Write

$$r_{avg}(x) = \int_{\mathbb{S}^{d-1}} r_\theta(\theta^*(x)) d\hat{\sigma}_K(\theta), \quad s_{avg}(y) = \int_{\mathbb{S}^{d-1}} s_\theta(\theta^*(y)) d\hat{\sigma}_K(\theta).$$

Given that $\varphi_1^\circ(f_{avg} + r_{avg}) = \varphi_1^\circ(f_{avg}) + r_{avg} \nabla \varphi_1^\circ(f_{avg}) + o(\|r_{avg}\|_\infty)$, the first order variation reads

$$\int r_{avg}(x) \nabla \varphi_1^\circ(f_{avg}(x) + \lambda^*(f_{avg}, g_{avg})) d\alpha(x) \quad (88)$$

$$+ \int s_{avg}(y) \nabla \varphi_2^\circ(g_{avg}(y) - \lambda^*(f_{avg}, g_{avg})) d\beta(y). \quad (89)$$

Then we define

$$\bar{\alpha} = \nabla \varphi_1^\circ(f_{avg} + \lambda^*(f_{avg}, g_{avg})) \alpha, \quad \bar{\beta} = \nabla \varphi_2^\circ(g_{avg} - \lambda^*(f_{avg}, g_{avg})) \beta,$$

such that the first order variation reads

$$\int r_{avg}(x) d\bar{\alpha}(x) + \int s_{avg}(y) d\bar{\beta}(y). \quad (90)$$

One can then explicit the definition of (r_{avg}, s_{avg}) , such that it reads

$$\int_{\mathbb{S}^{d-1}} \int r_\theta(\theta^*(x)) d\bar{\alpha}(x) + \int_{\mathbb{S}^{d-1}} \int s_\theta(\theta^*(y)) d\bar{\beta}(y) \quad (91)$$

$$= \int_{\mathbb{S}^{d-1}} \int r_\theta d\theta_\#^* \bar{\alpha}(x) + \int_{\mathbb{S}^{d-1}} \int s_\theta d\theta_\#^* \bar{\beta}(y). \quad (92)$$

By optimizing the above over the constraint set $\{r_\theta \oplus s_\theta \leq C_1\}$, we identify the computation of $\text{SOT}(\bar{\alpha}, \bar{\beta})$, which concludes the proof.

□

B.2 Frank-Wolfe methodology for computing UOT

Background: FW for UOT. Our approach to compute SUOT and USOT takes inspiration from the construction of (Séjourné et al., 2022b). It consists in applying a Frank-Wolfe (FW) procedure over the dual formulation of UOT. Such approach is equivalent to solve a sequence of balanced OT problems between measures $(\tilde{\alpha}, \tilde{\beta})$ which are iterative renormalizations of (α, β) . While the idea holds in wide generality, it is especially efficient in 1D where OT has low algorithmic complexity, and we reuse it in our sliced setting.

FW algorithm consists in optimizing a functional \mathcal{H} over a compact, convex set \mathcal{C} by optimizing its linearization $\nabla\mathcal{H}$. Given a current iterate x^t of FW algorithm, one computes $r^{t+1} \in \arg \max_{r \in \mathcal{C}} \langle \nabla\mathcal{H}(x^t), r \rangle$, and performs a convex update $x^{t+1} = (1 - \gamma_{t+1})x^t + \gamma_{t+1}r^{t+1}$. One typically chooses the learning rate $\gamma_t = \frac{2}{2+t}$. This yields the routine `FWStep` of Section 4 which is detailed below.

Algorithm 4 – FWStep(f, g, r, s, γ)

Input: $\alpha, \beta, f, g, \gamma$

Output: Normalized measures (α, β) as in (96)

```

 $f(x) \leftarrow (1 - \gamma)f(x) + \gamma r(x)$ 
 $g(y) \leftarrow (1 - \gamma)g(y) + \gamma s(y)$ 
Return  $(f, g)$ 

```

In the setting of UOT, one would take $\mathcal{C} = \{f \oplus g \leq C_d\}$. However, this set is not compact as it contains $(\lambda, -\lambda)$ for any $\lambda \in \mathbb{R}$. Thus, Séjourné et al. (2022b) propose to optimise a *translation-invariant* dual functional $\mathcal{H}(f, g; \alpha, \beta) \triangleq \sup_{\lambda \in \mathbb{R}} \mathcal{D}(f + \lambda, g - \lambda; \alpha, \beta)$, with \mathcal{D} defined in (3). Similar to the balanced OT dual, one has $\mathcal{H}(f + \lambda, g - \lambda; \alpha, \beta) = \mathcal{H}(f, g; \alpha, \beta)$, thus one can apply (Santambrogio, 2015, Proposition 1.11) to assume w.l.o.g. that, e.g., $f(0) = 0$ and restrict to a compact set of functions. We emphasize that FW algorithm is well-posed to optimize \mathcal{H} , but not \mathcal{D} .

Note that once we have the dual variables (f, g) maximizing \mathcal{H} , we retrieve optimal dual variables maximizing \mathcal{D} as $(f + \lambda^*(f, g), g - \lambda^*(f, g))$, where $\lambda^*(f, g) \triangleq \arg \max_{\lambda \in \mathbb{R}} \mathcal{D}(f + \lambda, g - \lambda; \alpha, \beta)$. The KL setting where $D_{\varphi_1} = \rho_1 \text{KL}$ and $D_{\varphi_2} = \rho_2 \text{KL}$ is especially convenient, because $\lambda^*(f, g)$ admits a closed form, which avoids iterative subroutines to compute it. In that case, it reads

$$\lambda^*(f, g) = \frac{\rho_1 \rho_2}{\rho_1 + \rho_2} \log \left(\frac{\int e^{-f(x)/\rho_1} d\alpha(x)}{\int e^{-g(y)/\rho_2} d\beta(y)} \right). \quad (93)$$

We summarize the FW algorithm for UOT in the proposition below. We refer to (Séjourné et al., 2022b) for more details on the algorithm and pseudo-code. We adapt this approach and result for SUOT and USOT.

Proposition B.1. (Séjourné et al., 2022b) *Assume φ° is smooth. Given current iterates $(f^{(t)}, g^{(t)})$, the linear FW oracle of UOT(α, β) is OT($\tilde{\alpha}^{(t)}, \tilde{\beta}^{(t)}$), where $\tilde{\alpha}^{(t)} = \nabla\varphi^\circ(f^{(t)} + \lambda^*(f^{(t)}, g^{(t)}))\alpha$ and $\tilde{\beta}^{(t)} = \nabla\varphi^\circ(g^{(t)} - \lambda^*(f^{(t)}, g^{(t)}))\beta$. In particular, one has $m(\tilde{\alpha}^{(t)}) = m(\tilde{\beta}^{(t)})$, thus the balanced OT problem always has finite value. More precisely, the FW update reads*

$$(f^{(t+1)}, g^{(t+1)}) = (1 - \gamma^{(t+1)})(f^{(t)}, g^{(t)}) + \gamma^{(t+1)}(r^{(t+1)}, s^{(t+1)}), \quad (94)$$

$$\text{where } (r^{(t+1)}, s^{(t+1)}) \in \arg \max_{r \oplus s \leq C_d} \int r(x) d\tilde{\alpha}^{(t)}(x) + \int s(y) d\tilde{\beta}^{(t)}(y). \quad (95)$$

Recall that the in KL setting one has $\varphi_i^\circ(x) = \rho_i(1 - e^{-x/\rho_i})$, thus $\nabla\varphi_i^\circ(x) = e^{-x/\rho_i}$. Thus in that case one normalizes the measures as

$$\tilde{\alpha} = \exp \left(-\frac{f + \lambda^*(f, g)}{\rho_1} \right) \alpha, \quad \tilde{\beta} = \exp \left(-\frac{g - \lambda^*(f, g)}{\rho_2} \right) \beta, \quad (96)$$

where λ^* is defined in (93).

This defines the `Norm` routine in Algorithm 1.

B.3 Implementation of Sliced OT to return dual potentials

Recall from Section 4, Algorithms 2 and 3 and more precisely, Propositions 4.1 and 4.2, that FW linear oracle is a sliced OT program, *i.e.*, a set of OT problems computed between univariate distributions of $\mathcal{M}_+(\mathbb{R})$. Therefore, a key building block of our algorithm is to compute the loss and dual variables of these univariate OT problems. We explain below how one can compute the sliced OT loss and dual potentials. The computation of the loss consists in implementing closed formulas of OT between univariate distributions, as detailed in (Santambrogio, 2015, Proposition 2.17). More precisely, when $C_1(x, y) = |x - y|^p$ and $(\mu, \nu) \in \mathcal{M}_+(\mathbb{R})$, then

$$\text{OT}(\mu, \nu) = \int_0^1 |F_\mu^{[-1]}(t) - F_\nu^{[-1]}(t)|^p dt, \quad (97)$$

where $F_\mu^{[-1]}$ denotes the inverse cumulative distribution function (ICDF) of μ .

Algorithm 5 – SlicedOTLoss($\alpha, \beta, \{\theta\}, p$)

Input: α, β , projections $\{\theta\}$, exponent p

Output: $\text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ as in (97)

for $\theta \in \{\theta\}$ **do**

 Project support of $\theta_\#^* \alpha$ and $\theta_\#^* \beta$

 Sort weights of $(\theta_\#^* \alpha, \theta_\#^* \beta)$ and support $(\theta^*(x), (\theta^*(y)))$ s.t. support is non-decreasing

 Compute ICDF of $\theta_\#^* \alpha$ and $\theta_\#^* \beta$

 Compute $\text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$ as in (97) with exponent p

end for

To compute dual potentials using backpropagation, one computes the sliced OT losses (using Algorithm 5) then calls the backpropagation w.r.t to inputs (α, β) , because their gradients are optimal dual potentials (Santambrogio, 2015, Proposition 7.17). We describe this procedure in Algorithm 6.

Algorithm 6 – SlicedOTPotentialsBackprop($\alpha, \beta, \{\theta\}, p$)

Input: α, β , projections $\{\theta\}$, exponent p

Output: Dual potentials (f_θ, g_θ) solving $\text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$

 Enable gradients w.r.t. $(\theta_\#^* \alpha, \theta_\#^* \beta)$

 Call $\text{SlicedOTLoss}(\alpha, \beta, \{\theta\}, p)$

 Sum (but do not average) losses $\mathcal{L} = \sum_\theta \text{OT}(\theta_\#^* \alpha, \theta_\#^* \beta)$.

 Backpropagate \mathcal{L} w.r.t. (α, β)

 Return (f_θ, g_θ) as gradients of \mathcal{L} w.r.t. (α, β) .

The implementation of the dual potentials using 1D closed forms relies on the north-west corner rule principle, which can be vectorized in PyTorch in order to be computed in parallel. The contribution of our implementation thus consists in making such algorithm GPU-compatible and allowing for a parallel computation for every slice simultaneously. We stress that this constitutes a non-trivial piece of code, and we refer the interested reader to the code in our supplementary material for more details on the implementation.

B.4 Output optimal sliced marginals

In all our algorithms, we focus on dual formulations of SUOT and USOT, which optimize the dual potentials. However, one might want the output variables of the primal formulation (See Definition 3.6). In particular, the marginals of optimal transport plans are interesting because they are interpreted as normalized versions of inputs (α, β) where geometric outliers have been removed. We detail where this interpretation comes from in the setting of UOT, and then give how it is adapted to SUOT and USOT. In particular, we justify that the `Norm` routine suffices to compute them.

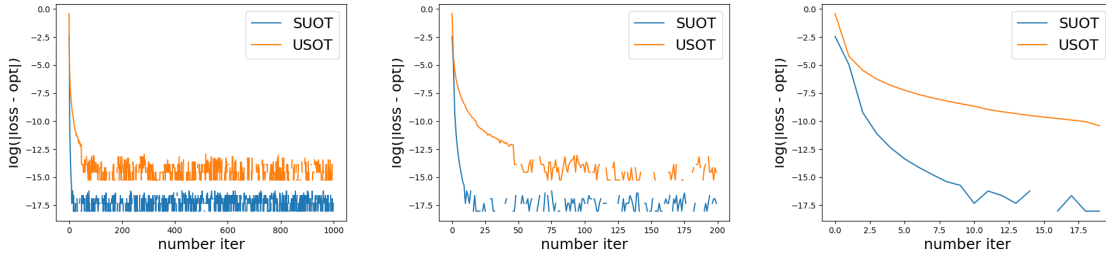


Figure 5: $|\text{SUOT}(\alpha, \beta) - \widehat{\text{SUOT}}_t|$ and $|\text{USOT}(\alpha, \beta) - \widehat{\text{USOT}}_t|$ against iteration t , where $\widehat{\text{SUOT}}_t, \widehat{\text{USOT}}_t$ are the estimated SUOT, USOT using t FW iterations. Plots are in log-scale. All figures are issued from the same run, but zoomed on a subset of first iterations: (*left*) 1000 iterations of FW, (*middle*) 200 iterations, (*right*) 20 iterations.

Case of UOT. We focus on the $D_{\varphi_i} = \rho_i \text{KL}$. As per (Liero et al., 2018, Equation 4.21), we have at optimality that the optimal transport π^* plan solving $\text{UOT}(\alpha, \beta)$ as in (2) has marginals (π_1^*, π_2^*) which read $\pi_1^* = e^{-f^*/\rho_1} \alpha$ and $\pi_2^* = e^{-g^*/\rho_2} \beta$, where (f^*, g^*) are the optimal dual potentials solving (3). Since on $\text{supp}(\pi^*)$ one also has $f^*(x) + g^*(y) = C_d(x, y)$, if the transportation cost $C_d(x, y)$ is large (*i.e.*, we are matching a geometric outlier), so are $f^*(x)$ and $g^*(y)$, and eventually the weights $\pi_1^*(x)$ and $\pi_2^*(y)$ are small, hence the interpretation of the geometric normalization of the measures. Note that in that case, one obtain (π_1^*, π_2^*) by calling $\text{Norm}(\alpha, \beta, f^*, g^*, \rho_1, \rho_2)$.

Case of SUOT. Since $\text{SUOT}(\alpha, \beta)$ consists in integrating $\text{UOT}(\theta_{\sharp}^* \alpha, \theta_{\sharp}^* \beta)$ w.r.t. σ , it shares many similarities with UOT. For any θ , we consider π_{θ} and (f_{θ}, g_{θ}) solving the primal and dual formulation of $\text{UOT}(\theta_{\sharp}^* \alpha, \theta_{\sharp}^* \beta)$. The marginals of π_{θ} are thus given by $(e^{-f_{\theta}/\rho_1} \alpha, e^{-g_{\theta}/\rho_2} \beta)$. In particular, we retrieve the observation made in Figure 1 that the optimal marginals change for each θ . In that case we call for each θ the routine $\text{Norm}(\alpha, \beta, f_{\theta}, g_{\theta}, \rho_1, \rho_2)$.

Case of USOT. Recall that the optimal marginals (π_1, π_2) in $\text{USOT}(\alpha, \beta)$ do not depend on θ , contrary to $\text{SUOT}(\alpha, \beta)$. Leveraging the dual formulation of Theorem 3.9, and looking at the Lagrangian which is defined in the proof of Theorem 3.9 (see Appendix A.2), we have the optimality condition that $\pi_1 = e^{-f_{\text{avg}}/\rho_1} \alpha$ and $\pi_2 = e^{-g_{\text{avg}}/\rho_2} \beta$. Thus in that case, calling $\text{Norm}(\alpha, \beta, f_{\text{avg}}, g_{\text{avg}}, \rho_1, \rho_2)$ yields the desired marginals.

B.5 Convergence of Frank-Wolfe iterations: Empirical analysis

We display below an experiment on synthetic dataset to illustrate the convergence of Frank-Wolfe iterations. We also provide insights on the number of iterations that yields a reasonable approximation: a few iterations suffices in our practical settings, typically $F = 20$.

The results are displayed in Figure 5. We consider the empirical distributions (α, β) computed over respectively, $N = 400$ and $M = 500$ samples over the unit hypercube $[0, 1]^d$, $d = 10$. Moreover, β is slightly shifted by a vector of uniform coordinates $0.5 \times \mathbf{1}_d$. We choose $\rho = 1$ and report the estimation of $\text{SUOT}(\alpha, \beta)$ and $\text{USOT}(\alpha, \beta)$ through Frank-Wolfe iterations. We estimate the true values by running $F = 5000$ iterations, and display the difference between the estimated score and the 'true' values. Appendix B.5 shows that numerical precision is reached in a few tens of iterations. As learning tasks do not usually require an estimation of losses up to numerical precision, we think that it is hence reasonable to take $F \approx 20$ in numerical applications.

Table 2: Dataset characteristics.

	BBCSport	Movies	Goodreads genre	Goodreads like
Doc	737	2000	1003	1003
Train	517	1500	752	752
Test	220	500	251	251
Classes	5	2	8	2
Mean words by doc	116 ± 54	182 ± 65	1491 ± 538	1491 ± 538
Median words by doc	104	175	1518	1518
Max words by doc	469	577	3499	3499

C Additional details on Section 5

C.1 Document classification: Technical details and additional results

C.1.1 Datasets

We sum up the statistics of the different datasets in Table 2.

BBCSport. The BBCSport dataset (Kusner et al., 2015) contains articles between 2004 and 2005, and is composed of 5 classes. We average over the 5 same train/test split of (Kusner et al., 2015). The dataset can be found in <https://github.com/mkusner/wmd/tree/master>.

Movie Reviews. The movie reviews dataset (Pang et al., 2002) is composed of 1000 positive and 1000 negative reviews. We take five different random 75/25 train/test split. The data can be found in <http://www.cs.cornell.edu/people/pabo/movie-review-data/>.

Goodreads. This dataset, proposed in (Maharjan et al., 2017), and which can be found at https://ritual.uh.edu/multi_task_book_success_2017/, is composed of 1003 books from 8 genres. A first possible classification task is to predict the genre. A second task is to predict the likability, which is a binary task where a book is said to have success if it has an average rating ≥ 3.5 on the website Goodreads (<https://www.goodreads.com>). The five train/test split are randomly drawn with 75/25 proportions.

C.1.2 Technical Details

All documents are embedded with the Word2Vec model (Mikolov et al., 2013) in dimension $d = 300$. The embedding can be found in <https://drive.google.com/file/d/0B7XkCwpI5KDYN1NUTTLSS21pQmM/view?resourcekey=0-wjGZdNAUop6WYkTtMip30g>.

In this experiment, we report the results averaged over 5 random train/test split. For discrepancies which are approximated using random projections, we additionally average the results over 3 different computations, and we report this standard deviation in Table 1. Furthermore, we always use 500 projections to approximate the sliced discrepancies. For Frank-Wolfe based methods, we use 10 iterations, which we found to be enough to have a good accuracy. We added an ablation of these two hyperparameters in Figure 7. We report the results obtained with the best ρ for USOT and SUOT computed among a grid $\rho \in \{10^{-4}, 5 \cdot 10^{-4}, 10^{-3}, 5 \cdot 10^{-3}, 10^{-2}, 10^{-1}, 1\}$. For USOT, the best ρ is consistently around $5 \cdot 10^{-3}$ for the Movies and Goodreads datasets, and around $5 \cdot 10^{-4}$ for the BBCSport dataset. We used a second finer grid and reported the results obtained with $\rho = 0.00021$ on BBCSport, $\rho = 0.004$ for Goodreads on the likability task and $\rho = 0.003$ for the genre task. For SUOT, the best ρ obtained was 0.01 for the BBCSport dataset, 1.0 for the movies dataset and 0.5 for the goodreads dataset. For UOT, we used $\rho = 1.0$ on the BBCSport dataset. For the movies dataset, the best ρ obtained on a subset was 50, but it took an unreasonable amount of time to run on the full dataset as the runtime increases with ρ (see (Chapel et al., 2021, Figure 3)). On the goodreads dataset, it took too much memory on the GPU. For Sinkhorn UOT, we used $\varepsilon = 0.1$ and $\rho = 1.0$ on the BBCSport and $\varepsilon = 0.001$, $\rho = 0.1$ on the Goodreads dataset, and $\varepsilon = 0.01$ and $\rho = 0.1$ on the Movies dataset. Note

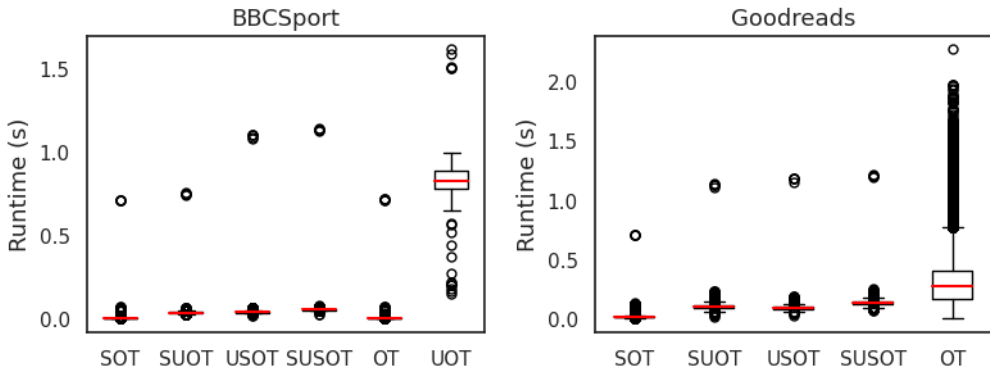


Figure 6: Runtime on the BBCSport dataset (left) and on the Goodreads dataset (right).

that we also tested other set of parameters for UOT and Sinkhorn UOT, but on a coarser grid than USOT and SUOT given their computational time. For each method, the number of neighbors used for the k-NN method is obtained via cross-validation.

C.1.3 Additional experiments

Runtime. We report in Figure 6 the runtime of computing the different discrepancies between each pair of documents, and in Table 3 the full runtimes. On the BBCSport dataset, the documents have in average 116 words, thus the main bottleneck is the projection step for sliced OT methods. Hence, we observe that OT runs slightly faster than SOT and the sliced unbalanced counterparts. Goodreads is a dataset with larger documents, with on average 1491 words by document. Therefore, as OT scales cubically with the number of samples, we observe here that all sliced methods run faster than OT, which confirms that sliced methods scale better w.r.t. the number of samples. In this setting, we were not able to compute UOT with the POT implementation in a reasonable time. Computations have been performed with a NVIDIA Tesla V100 GPU.

Table 3: Runtimes on Document Classification

		BBCSport	Goodreads
OT	Average ($\cdot 10^{-3}$ s)	$3.29_{\pm 1.61}$	$440.30_{\pm 259}$
	Full (s)	891	221252
SOT	Average ($\cdot 10^{-3}$ s)	$1.80_{\pm 6.22}$	$4.49_{\pm 1.44}$
	Full (s)	487	2256
USOT	Average ($\cdot 10^{-3}$ s)	$14.67_{\pm 1.29}$	$14.45_{\pm 0.88}$
	Full (s)	3897	7260
SUOT	Average ($\cdot 10^{-3}$ s)	$13.9_{\pm 1.21}$	$14.32_{\pm 0.95}$
	Full (s)	3770	7193

Ablations. We plot in Figure 7 accuracy as a function of the number of projections and the number of iterations of the Frank-Wolfe algorithm. We averaged the accuracy obtained with the same setting described in Appendix C.1.2, with varying number of projections $K \in \{4, 10, 21, 46, 100, 215, 464, 1000\}$ and number of FW iterations $F \in \{1, 2, 3, 4, 5, 10, 15, 20\}$. Regarding the hyperparameter ρ , we selected the one returning the best accuracy, *i.e.*, $\rho = 5 \cdot 10^{-4}$ for USOT and $\rho = 10^{-2}$ for SUOT.

Choice of ρ . In Table 4, we also add the results obtained when ρ is chosen by cross validation for USOT and SUOT. In this case, the results are slightly below the best one, but are still better than SOT.

Unnormalizing measures. As we have mentioned in Section 5, we have performed document classification experiments by normalizing word histograms to be probabilities. It allows to compare SUOT and USOT with SOT since sliced OT is only well-defined between probabilities. However, it seems reasonable to compare

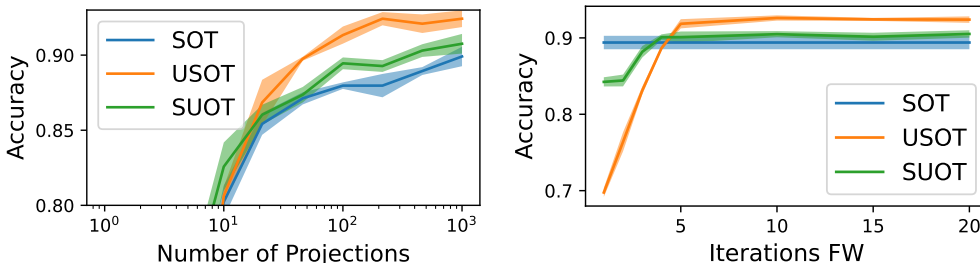


Figure 7: Ablation on BBCSport of the number of projections (*left*) and of the number of Frank-Wolfe iterations (*right*).

Table 4: Accuracy on document classification. Grid-search over ρ is performed. We report the best accuracy over ρ , which corresponds to $\rho = 5.10^{-6}$ for Unnormalized USOT and $\rho = 5.10^{-6}$ for SOPT.

	BBCSport	Movies	Goodreads genre	Goodreads like
OT	94.55	74.44	55.22	71.00
UOT	96.73	-	-	-
Sinkhorn UOT	95.45	72.48	53.55	67.81
SOT	89.39 \pm 0.76	66.95 \pm 0.45	50.09 \pm 0.51	65.60 \pm 0.20
SUOT	90.12 \pm 0.15	67.84 \pm 0.37	50.15 \pm 0.04	66.72 \pm 0.38
USOT	93.52 \pm 0.04	69.21 \pm 0.37	52.67 \pm 0.62	67.78 \pm 0.39
SUSOT	92.73 \pm 0.27	69.53 \pm 0.53	51.93 \pm 0.53	67.33 \pm 0.26
SUOT (+CV on ρ)	90.00 \pm 0.59	67.40 \pm 0.64	49.67 \pm 0.79	66.43 \pm 0.44
USOT (+CV on ρ)	92.61 \pm 0.55	68.64 \pm 0.29	52.06 \pm 7.20	66.61 \pm 0.72
USOT (Unnormalized)	86 \pm 0.56	-	-	-
SOPT (Unnormalized)	87.27 \pm 0.20	-	-	-

with other unbalanced sliced methods such as SOPT (Bai et al., 2023). We chose to compare with this competitor since their code is available in Python. However, a numerical restriction of their algorithm is that it only outputs measures with constants weights, *i.e.*, distributions $\alpha = \sum \alpha_i \delta_{x_i}$ and $\beta = \sum \beta_j \delta_{y_j}$ where $\alpha_i = \beta_j = 1$, but the number of samples in α and β may differ. Under this modeling assumption, the total mass of each measure corresponds to the number of words in the sentence. We performed the comparison on the BBC dataset, using 500 projections for both SOPT and USOT. Unfortunately, the quadratic footprint of computing the similarity kernel does not scale reasonably for SOPT for larger datasets such as Movies or Goodreads, especially because their algorithm is not GPU-compatible compared to ours. We cross-validated the parameter $\rho \in \{p.10^{-k}, k \in \llbracket 0, 6 \rrbracket, p \in \{1., 5.\}\}$

The result is detailed in the table below. What is noticeable is that the performance degrades for both USOT and SOPT using this parametrization. Furthermore, we observed that the parameter ρ yielding the best accuracy is much smaller for unnormalized measures than for the best one for normalized histograms (*i.e.*, $1e - 5$ here compared to $1e - 3$ with normalized measures). Our interpretation of this observation is that considering unnormalized measures adds an additional information of the sentence length via the masses of (α, β) . It seems that this additional information dominates the comparison of measures, instead of focusing on the measures support (*i.e.*, the word embedding) which encodes the semantic information of words. When ρ is large the kernel value of USOT/SOPT is mainly dictated by the mass (*i.e.*, sentence length) comparison. Thus smaller ρ seems to give less importance on sentence length, hence a better performance. We also note that performance of SOPT and USOT on unnormalized measures are rather similar. It means that for the choice of marginal prior $D_\varphi = \rho TV$ or $D_\varphi = \rho KL$ does not significantly matter for this specific task, compared to the preprocessing normalization of measures.

C.2 Unbalanced sliced Wasserstein barycenters

We define below the formulation of the USOT barycenter which was used in the experiments of Figure 4 to average predictions of geophysical data. We then detail how we computed it.

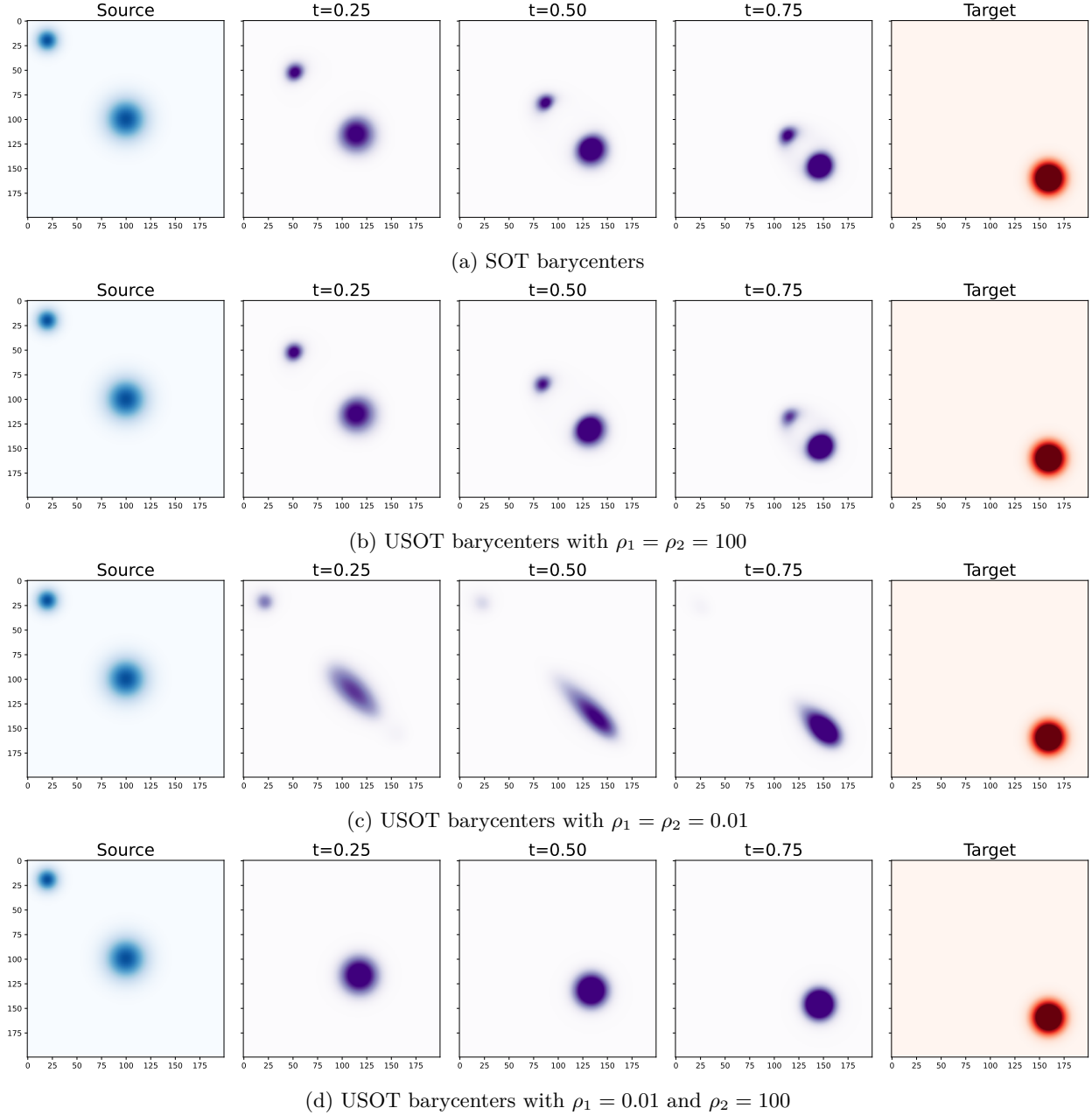


Figure 8: **Interpolation with USOT as a barycenter computation.** We compare different interpolations using SOT or USOT with different settings for the ρ values

Definition C.1. Consider a set of measures $(\alpha_1, \dots, \alpha_B) \in \mathcal{M}_+(\mathbb{R}^d)^B$, and a set of non-negative coefficients $(\omega_1, \dots, \omega_B) \geq 0$ such that $\sum_{b=1}^B \omega_b = 1$. We define the barycenter problem (in the KL setting) as

$$\mathcal{B}((\alpha_b)_b, (\omega_b)_b) \triangleq \inf_{\beta \in \mathcal{P}(\mathbb{R}^d)} \sum_{b=1}^B \omega_b \text{USOT}(\alpha_b, \beta), \quad (98)$$

$$= \inf_{\beta \in \mathcal{P}(\mathbb{R}^d)} \sum_{b=1}^B \inf_{(\pi_{b,1}, \pi_{b,2})} \text{SOT}(\pi_{b,1}, \pi_{b,2}) + \rho_1 \text{KL}(\pi_{b,1} | \alpha_b) + \rho_2 \text{KL}(\pi_{b,2} | \beta), \quad (99)$$

where $\mathcal{P}(\mathbb{R}^d)$ denotes the set of probability measures.

To compute the barycenter, we aggregate several building blocks. First, since we consider that the barycenter $\beta \in \mathcal{P}(\mathbb{R}^d)$ is a probability, we perform mirror descent as in (Beck & Teboulle, 2003; Cuturi & Doucet, 2014b). More precisely, we use a Nesterov accelerated version of mirror descent. We also tried projected gradient descent, but it did not yield consistent outputs (due to convergence speed (Beck & Teboulle, 2003)). Second, we use a Stochastic-USOT version (see Section 4), *i.e.*, we sample new projections at each iteration of the barycenter update (but not a each iteration of the FW subroutines in Algorithm 3). This procedure is described in Algorithm 7.

Algorithm 7 – Barycenter $((\alpha_b)_b, (\omega_b)_b, \rho_1, \rho_2, lr)$

Input: measures $(\alpha_b)_b$, weights $(\omega_b)_b$, ρ_1, ρ_2 , learning rate lr , FW iter F

Output: Optimal barycenter β of (98)

```

 $t \leftarrow 1$ 
Init  $(\beta, \tilde{\beta}, \hat{\beta})$  as uniform distribution over a grid
while not converged do do
   $\gamma \leftarrow \frac{2}{(t+1)}$ ,
   $\beta \rightarrow (1 - \gamma)\hat{\beta} + \gamma\tilde{\beta}$ 
  Sample projections  $(\theta_k)_{k=1}^K$ 
  Compute  $\mathcal{B}((\alpha_b)_b, (\omega_b)_b)$  by calling USOT $(\alpha_b, \beta, F, (\theta_k)_{k=1}^K, \rho_1, \rho_2)$  in Algorithm 3 for each  $b$ 
  Compute  $g$  as the gradient of  $\mathcal{B}((\alpha_b)_b, (\omega_b)_b)$  w.r.t. variable  $\beta$ 
   $\tilde{\beta} \leftarrow \exp(-lr \times \gamma^{-1} \times g)\beta$ 
   $\tilde{\beta} \leftarrow \tilde{\beta}/m(\tilde{\beta})$ 
   $\hat{\beta} \leftarrow (1 - \gamma)\hat{\beta} + \gamma\tilde{\beta}$ 
   $t \leftarrow t + 1$ 
end while

```

We illustrate this algorithm with several examples of interpolation in Figure 8. We propose to compute an interpolation between two measures located on a fixed grid of size 200×200 with different values of ρ_i in $D_{\varphi_i} = \rho_i \text{KL}$. For illustration purposes, we construct the *source* distribution as a mixture of two Gaussians with a small and a larger mode, and the *target* distribution as a single Gaussian. Those distributions are normalized over the grid such that both total norms are equal to one (which is not required by our unbalanced sliced variants but grants more interpretability and possible comparisons with SOT). Figure 8a shows the result of the interpolation at three timestamps ($t = 0.25, 0.5$ and 0.75) of a SOT interpolation (within this setting, $\omega_1 = 1 - t$ and $\omega_2 = t$). As expected, the two modes of the source distribution are transported over the target one. We verify in Figure 8b that for a large value of $\rho_1 = \rho_2 = 100$, the USOT interpolation behaves similarly as SOT, as expected from the theory. When $\rho_1 = \rho_2 = 0.01$, the smaller mode is not moved during the interpolation, whereas the larger one is stretched toward the target (Figure 8c). Finally, in Figure 8d, an asymmetric configuration of $\rho_1 = 0.01$ and $\rho_2 = 100$ allows to get an interpolation when only the big mode of the source distribution is displaced toward the target. In all those cases, the mirror-descent algorithm 7 is run for 500 iterations. Even for a large grid of 200×200 , those different results are obtained in a 2 – 3 minutes on a commodity GPU, while the OT or UOT barycenters are untractable with a limited computational budget.

C.3 Unbalanced version of hyperbolic SOT.

To illustrate the modularity of our FW algorithm, we aim at comparing synthetic mixtures of Wrapped Normal Distribution on the 2-hyperbolic manifold \mathbb{H} (Nagano et al., 2019), so that the FW oracle is hyperbolic sliced OT (Bonet et al., 2023c). The parameter θ characterizes on \mathbb{H} any geodesic curve passing through the origin, and each sample is projected by taking the shortest path to such geodesics. Once projected on a geodesic curve, we sort data and compute SOT w.r.t. hyperbolic metric $d_{\mathbb{H}}$. We consider the 2-hyperbolic manifold on the Poincaré disc. As illustrated in Figure 9, the input measure α (in red) is a mixture of 3 isotropic normal distributions, with a mode at the top of the disc playing the role of an outlier. The measure β is a mixture of two anisotropic normal distributions, whose means are close to two modes of α , but are slightly shifted at the disc’s center. We show on Figure 9 the impact of the parameter $\rho = \rho_1 = \rho_2$ on the optimal marginals of USOT.

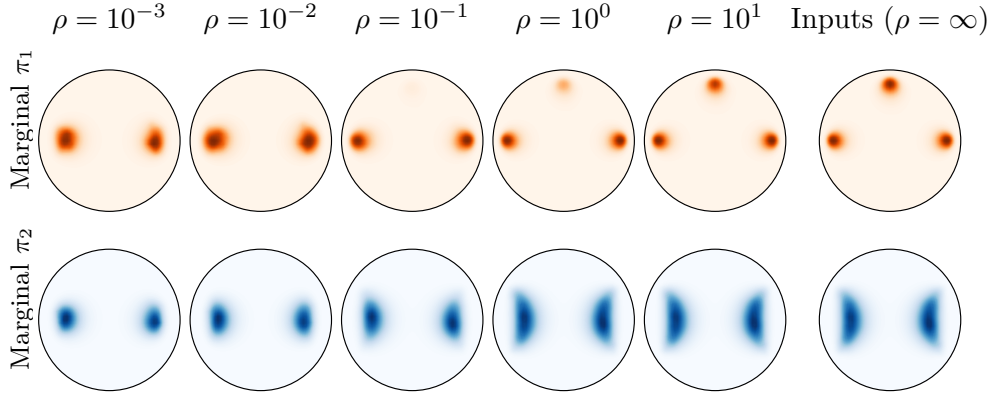


Figure 9: KDE estimation (kernel $e^{-d_{\mathbb{H}}^2/\sigma}$) of optimal (π_1, π_2) of $\text{USOT}(\alpha, \beta)$ when $D_{\varphi_i} = \rho\text{KL}$.

This experiment illustrates several take-home messages, mentioned in Section 3. First, the optimal marginals (π_1, π_2) are renormalisation of (α, β) accounting for their geometry, which are able to remove outliers for properly tuned ρ . When ρ is large, $(\pi_1, \pi_2) \simeq (\alpha, \beta)$ and we retrieve SOT. When ρ is too small, outliers are removed, but we see a shift of the modes, so that modes of (π_1, π_2) are closer to each other, but do not exactly correspond to those of (α, β) . Second, note that such plot cannot be made with SUOT, since the optimal marginals depend on the projection θ (see Figure 1). Third, we emphasize that we are indeed able to reuse any variant of SOT.

C.4 Choice and interpretation of hyperparameter ρ

An immediate drawback of our framework is the induced additional computational cost w.r.t. SOT. While the above experimental results show that SUOT and USOT improve performance significantly over SOT, and though the complexity is still sub-quadratic in number of samples, our FW approach uses SOT as a subroutine, rendering it necessarily more expensive. Additionally, another practical burden comes from the introduction of extra hyperparameters (ρ_1, ρ_2) , which may be tuned using cross-validation. Therefore, a future direction would be to derive efficient strategies to tune (ρ_1, ρ_2) , maybe w.r.t. the applicative context, and further complement the possible interpretations of ρ as a 'threshold' for the geometric information encode by C_1, C_d . While we leave the automation of tuning (ρ_1, ρ_2) for future works, we provide below some details and intuitions on the choice of ρ for the previous experiments. We hope these insights will help the practitioner on how they should chose tune this additional parameter.

General intuition on ρ . The parameter ρ when $\rho_1 = \rho_2 = \rho$ can be understood as a *characteristic distance* to decide whether or not two sample should be matched by the coupling π in the primal formulation of (2). Typically, transportation happens for samples (x, y) such that $C_d(x, y) \leq \rho$, while samples such that $C_d(x, y) \geq \rho$ are interpreted as geometric outliers, and are discarded in the matching $\pi(x, y)$. In the case of SUOT and USOT, there is somehow a similar interpretation, but not for the same quantities, and we rely on their definitions (Equations 7 and 10), as well as the constraint set \mathcal{E} in Theorem 3.9.

One sees that for $\text{SUOT}(\alpha, \beta)$ we have a set of 1D-UOT problems between $(\theta_{\#}^* \alpha, \theta_{\#}^* \beta)$, thus the threshold interpretation holds depending on whether $C_1(\theta^*(x), \theta^*(y)) \leq \rho$ or $C_1(\theta^*(x), \theta^*(y)) \geq \rho$. In particular the dependence in θ explains why the outlier threshold depends on the considered projection. Note also we consider C_1 instead of C_d .

For $\text{USOT}(\alpha, \beta)$ it is different because the marginals (π_1, π_2) which we optimize in Equation (9) are independent of θ , and common to all projections. Informally speaking, we interpret that the threshold value to discard a matching between (x, y) depends on whether some quantity proportional to $\int_{\mathbb{S}^{d-1}} C_1(\theta^*(x), \theta^*(y)) d\pi_{\theta}(x, y) d\sigma(\theta)$ is larger or smaller than ρ . This quantity is not properly defined as it depends on the optimized variables $(\pi_{\theta})_{\theta}$, hence the informality of our intuition. However, we wish to emphasize that the parameter ρ should be interpreted differently between SUOT and USOT. As highlighted

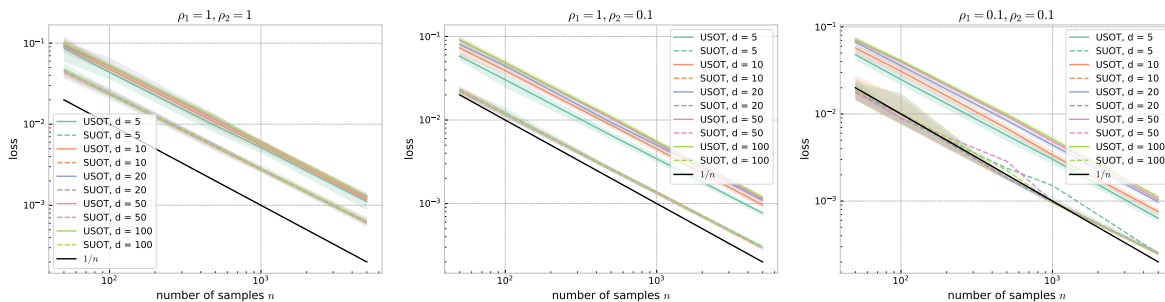


Figure 10: **Sample complexity:** $\mathbb{E}[|\mathcal{L}(\hat{\alpha}_n, \hat{\beta}_n) - \mathcal{L}(\alpha, \beta)|]$ against n , with $\mathcal{L} = \text{SUOT}$ or USOT and $\alpha = \beta = \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ (thus $\mathcal{L}(\alpha, \beta) = 0$). Results are averaged over 20 runs and the shaded areas represent the 10th and 90th percentiles. All curves exhibit the same convergence rate for any d , that is $1/n$ (see solid line black curve).

experimentally for document classification in Figure 2, we observe that the value of ρ yielding the best performance is not the same for each loss.

Choice of ρ for hyperbolic data. In Figure 9, the hyperbolic distance between overlapping modes is 0.96, while distance from side modes to the top red outlier is 2.83. Thus, a proper choice of should lie in between, which seems consistent with the observation of Figure 9 for $\rho = 1$. Indeed we see that we have a satisfying trade-off between removing the top mode and preserving the crescent shape structure of main blue modes.

Choice of ρ for barycenter experiments. For the barycenter, we used insights from Figure 8 which interpolates circular blobs using asymmetric (ρ_1, ρ_2) , where ρ_1 is the parameter penalizing the input measures fidelity, and ρ_2 the parameter of the barycenter. For Figure 4 (especially line (d)), we also took assymetric (ρ_1, ρ_2) with large $\rho_2 = 1e4$ for the barycenter to force data matching. Then for inputs $\rho_1 = 1e1$ is roughly the distance between cyclones (see Figure 4), to keep them in the barycenter. All in all, we force the barycenter to match the cyclone structure which matters most, while any structure who would be beyond this ρ_1 distance between input measure would be discarded.

Interpretation of ρ for document classification. In this task the measures' support are given by word embedding in high dimension, for which we have no intuition of what is for instance the characteristic distance between different semantic clusters, and thus no idea on how ρ should be tuned. For this reason (and more generally in ML tasks), we need to perform a cross-validation over this hyperparameter. We would like to comment the dependence of the document classification accuracy w.r.t. ρ , which can be observed in Figure 2. One can notice that as ρ increases, the accuracy increases until it reaches a 'peak', until then it decreases to reach a plateau as $\rho \rightarrow \infty$. When $\rho \rightarrow \infty$, SUOT and USOT converge to SOT (see Definitions 2.4 and 3.6), and we get similar performances. As $\rho \rightarrow 0$, marginals (π_1, π_2) are allowed to differ significantly from inputs (α, β) , meaning that SUOT/USOT almost ignore input data. Therefore, ρ should be tuned to extract information from inputs while removing noise. In Figure 2, the 'peaks' correspond to such optimal ρ , and the gain in performance justify the use of SUOT/USOT over SOT.

C.5 Illustration of the sample complexity

We investigate the sample complexity of SUOT and USOT in practice and report the results in Figure 10. Our goal is to empirically verify Theorem 3.4 for SUOT, and explore the convergence rate for USOT. To this end, we consider $\alpha = \beta = \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ and compute $\text{SUOT}(\alpha_n, \beta_n)$ and $\text{USOT}(\alpha_n, \beta_n)$ for different number of samples n and dimension d . This allows us to explore the convergence rate of $\text{SUOT}(\alpha_n, \beta_n)$ to $\text{SUOT}(\alpha, \beta) = 0$ (respectively, of $\text{USOT}(\alpha_n, \beta_n)$ to $\text{USOT}(\alpha, \beta) = 0$) as a function of n and d .

Figure 10 shows that all curves share the same slope w.r.t n , for any d and for both SUOT and USOT. This experiment is consistent with the dimension-free rate we established in Theorem 3.4 for SUOT. Interestingly, it also reveals that the dimension-free rate holds for USOT as well in that specific setting. More experiments and/or theoretical justification are needed to verify if this holds for more general distributions.