

---

# Equitable and Optimal Transport with Multiple Agents

---

**Meyer Scetbon\***

CREST, ENSAE,  
Inst. Polytechnique de Paris

**Laurent Meunier\***

Facebook AI Research  
Univ. Paris-Dauphine

**Jamal Atif**

CNRS, LAMSADE,  
Univ. Paris-Dauphine

**Marco Cuturi**

Google Brain  
CREST, ENSAE

## Abstract

We introduce an extension of the Optimal Transport problem when multiple costs are involved. Considering each cost as an agent, we aim to share equally between agents the work of transporting one distribution to another. To do so, we minimize the transportation cost of the agent who works the most. Another point of view is when the goal is to partition equitably goods between agents according to their heterogeneous preferences. Here we aim to maximize the utility of the least advantaged agent. This is a fair division problem. Like Optimal Transport, the problem can be cast as a linear optimization problem. When there is only one agent, we recover the Optimal Transport problem. When two agents are considered, we are able to recover Integral Probability Metrics defined by  $\alpha$ -Hölder functions, which include the widely-known Dudley metric. To the best of our knowledge, this is the first time a link is given between the Dudley metric and Optimal Transport. We provide an entropic regularization of that problem which leads to an alternative algorithm faster than the standard linear program.

## 1 Introduction

Optimal Transport (OT) has gained interest last years in machine learning with diverse applications in neuroimaging (Janati et al., 2020), generative models (Arjovsky et al., 2017; Salimans et al., 2018), supervised learning (Courty et al., 2016), word embeddings (Alvarez-Melis et al., 2018), reconstruction cell trajectories (Yang et al., 2020; Schiebinger et al., 2019) or

adversarial examples (Wong et al., 2019). The key to use OT in these applications lies in the gain of computation efficiency thanks to regularizations that smoothes the OT problem. More specifically, when one uses an entropic penalty, one recovers the so called Sinkhorn distances (Cuturi, 2013). In this paper, we introduce a new family of variational problems extending the optimal transport problem when multiple costs are involved with various applications in fair division of goods/work and operations research problems.

Fair division (Steinhaus, 1949) has been widely studied by the artificial intelligence (Lattimore et al., 2015) and economics (Moulin, 2004) communities. Fair division consists in partitioning diverse resources among agents according to some fairness criteria. One of the standard problems in fair division is the fair cake-cutting problem (Dubins and Spanier, 1961; Brandt et al., 2016). The cake is an heterogeneous resource, such as a cake with different toppings, and the agents have heterogeneous preferences over different parts of the cake, i.e., some people prefer the chocolate toppings, some prefer the cherries, others just want a piece as large as possible. Hence, taking into account these preferences, one might share the cake equitably between the agents. A generalization of this problem, for which achieving fairness constraints is more challenging, is when the splitting involves several heterogeneous cakes, and where the agents have linked preferences over the different parts of the cakes. This problem has many variants such as the cake-cutting with two cakes (Cloutier et al., 2010), or the Multi Type Resource Allocation (Mackin and Xia, 2015; Wang et al., 2019). In all these models it is assumed that there is only one indivisible unit per type of resource available in each cake, and once an agent choose it, he or she has to take it all. In this setting, the cake can be seen as a set where each element of the set represents a type of resource, for instance each element of the cake represents a topping. A natural relaxation of these problems is when a divisible quantity of each type of resources is available. We introduce EOT (**E**quitable and **O**ptimal **T**ransport), a formulation that solves both the cake-cutting and the cake-cutting with two cakes problems in this setting.

---

\* Authors contributed equally. Proceedings of the 24<sup>th</sup> International Conference on Artificial Intelligence and Statistics (AISTATS) 2021, San Diego, California, USA. PMLR: Volume 130. Copyright 2021 by the author(s).

Our problem expresses as an optimal transportation problem. Hence, we prove duality results and provide fast computation based on Sinkhorn algorithm. As interesting properties, some Integral Probability Metrics (IPMs) (Müller, 1997) as Dudley metric (Dudley et al., 1966), or standard Wasserstein metric (Villani, 2003) are particular cases of the EOT problem.

**Contributions.** In this paper we introduce EOT an extension of Optimal Transport which aims at finding an equitable and optimal transportation strategy between multiple agents. We make the following contributions:

- In Section 3, we introduce the problem and show that it solves a fair division problem where heterogeneous resources have to be shared among multiple agents. We derive its dual and prove strong duality results. As a by-product, we show that EOT is related to some usual IPMs families and in particular the widely known Dudley metric.
- In Section 4, we propose an entropic regularized version of the problem, derive its dual formulation, obtain strong duality. We then provide an efficient algorithm to compute EOT. Finally we propose other applications of EOT for Operations Research problems.

## 2 Related Work

**Optimal Transport.** Optimal transport aims to move a distribution towards another at lowest cost. More formally, if  $c$  is a cost function on the ground space  $\mathcal{X} \times \mathcal{Y}$ , then the relaxed Kantorovich formulation of OT is defined for  $\mu$  and  $\nu$  two distributions as

$$W_c(\mu, \nu) := \inf_{\gamma} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\gamma(x, y)$$

where the infimum is taken over all distributions  $\gamma$  with marginals  $\mu$  and  $\nu$ . Kantorovich theorem states the following strong duality result under mild assumptions (Villani, 2003)

$$W_c(\mu, \nu) = \sup_{f, g} \int_{\mathcal{X}} f(x) d\mu(x) + \int_{\mathcal{Y}} g(y) d\nu(y)$$

where the supremum is taken over continuous bounded functions satisfying for all  $x, y$ ,  $f(x) + g(y) \leq c(x, y)$ . The question of considering an optimal transport problem when multiple costs are involved has already been raised in recent works. For instance, (Paty and Cuturi, 2019) proposed a robust Wasserstein distance where the distributions are projected on a  $k$ -dimensional subspace that maximizes their transport cost. In that sense, they aim to choose the most expensive cost among Mahalanobis square distances with kernels of rank  $k$ . In

articles (Li et al., 2019; Sun et al., 2020), the authors aim to learn a cost given observed matchings by inverting the optimal transport problem (Dupuy et al., 2016). In (Petrovich et al., 2020) the authors study “feature-robust” optimal transport, which can be also seen as a robust cost selection for optimal transport. In articles (Genevay et al., 2017; Scetbon and Cuturi, 2020), the authors learn an adversarial cost to train a generative adversarial network. Here, we do not aim to consider a worst case scenario among the available costs but rather consider that the costs work together in order to split equitably the transportation problem among them at lowest cost.

**Entropic relaxation of OT.** Computing exactly the optimal transport cost requires solving a linear program with a supercubic complexity ( $n^3 \log n$ ) (Tarjan, 1997) that results in an output that is *not* differentiable with respect to the measures’ locations or weights (Bertsimas and Tsitsiklis, 1997). Moreover, OT suffers from the curse of dimensionality (Dudley, 1969; Fournier and Guillin, 2015) and is therefore likely to be meaningless when used on samples from high-dimensional densities. Following the line of work introduced by Cuturi (2013), we propose an approximated computation of our problem by regularizing it with an entropic term. Such regularization in OT accelerates the computation, makes the problem differentiable with regards to the distributions (Feydy et al., 2018) and reduces the curse of dimensionality (Genevay et al., 2018). Taking the dual of the approximation, we obtain a smooth and convex optimization problem under a simplicial constraint.

**Fair Division.** Fair division of goods has a long standing history in economics and computational choice. A classical problem is the fair cake-cutting that consists in splitting the cake between  $N$  individuals according to their heterogeneous preferences. The cake  $\mathcal{X}$ , viewed as a set, is divided in  $\mathcal{X}_1, \dots, \mathcal{X}_N$  disjoint sets among the  $N$  individuals. The utility for a single individual  $i$  for a slice  $S$  is denoted  $V_i(S)$ . It is often assumed that  $V_i(\mathcal{X}) = 1$  and that  $V_i$  is additive for disjoint sets. There exists many criteria to assess fairness for a partition  $\mathcal{X}_1, \dots, \mathcal{X}_N$  such as proportionality ( $V_i(\mathcal{X}_i) \geq 1/N$ ), envy-freeness ( $V_i(\mathcal{X}_i) \geq V_i(\mathcal{X}_j)$ ) or equitability ( $V_i(\mathcal{X}_i) = V_j(\mathcal{X}_j)$ ). The cake-cutting problem has applications in many fields such as dividing land estates, advertisement space or broadcast time. An extension of the cake-cutting problem is the cake-cutting with two cakes problem (Cloutier et al., 2010) where two heterogeneous cakes are involved. In this problem, preferences of the agents can be coupled over the two cakes. The slice of one cake that an agent prefers might be influenced by the slice of the other cake that he or

she might also obtain. The goal is to find a partition of the cakes that satisfies fairness conditions for the agents sharing the cakes. Cloutier et al. (2010) studied the envy-freeness partitioning. Both the cake-cutting and the cake-cutting with two cakes problems assume that there is only one indivisible unit of supply per element  $x \in \mathcal{X}$  of the cake(s). Therefore sharing the cake(s) consists in obtaining a partition of the set(s). In this paper, we show that EOT is a relaxation of the cutting cake and the cake-cutting with two cakes problems, when there is a divisible amount of each element of the cake(s). In that case, cakes are no more sets but distributions that we aim to divide between the agents according to their coupled preferences.

**Integral Probability Metrics.** In our work, we make links with some integral probability metrics. IPMs are (semi-)metrics on the space of probability measures. For a set of functions  $\mathcal{F}$  and two probability distributions  $\mu$  and  $\nu$ , they are defined as

$$\text{IPM}_{\mathcal{F}}(\mu, \nu) = \sup_{f \in \mathcal{F}} \int f d\mu - \int f d\nu.$$

For instance, when  $\mathcal{F}$  is chosen to be the set of bounded functions with uniform norm less or equal than 1, we recover the Total Variation distance (Steineman, 1983) (TV). They recently regained interest in the Machine Learning community thanks to their application to Generative Adversarial Networks (GANs) (Goodfellow et al., 2014) where IPMs are natural metrics for the discriminator (Dziugaite et al., 2015; Arjovsky et al., 2017; Mroueh and Sercu, 2017; Husain et al., 2019). They also helped to build consistent two-sample tests (Gretton et al., 2012; Scetbon and Varoquaux, 2019). However when a closed form of the IPM is not available, exact computation of IPMs between discrete distributions may not be possible or can be costly. For instance, the Dudley metric can be written as a Linear Program (Sriperumbudur et al., 2012) which has at least the same complexity as standard OT. Here, we show that the Dudley metric is in fact a particular case of our problem and obtain a faster approximation thanks to the entropic regularization.

### 3 Equitable and Optimal Transport

**Notations.** Let  $\mathcal{Z}$  be a Polish space, we denote  $\mathcal{M}(\mathcal{Z})$  the set of Radon measures on  $\mathcal{Z}$ . We call  $\mathcal{M}_+(\mathcal{Z})$  the sets of positive Radon measures, and  $\mathcal{M}_+^1(\mathcal{Z})$  the set of probability measures. We denote  $\mathcal{C}^b(\mathcal{Z})$  the vector space of bounded continuous functions on  $\mathcal{Z}$ . Let  $\mathcal{X}$  and  $\mathcal{Y}$  be two Polish spaces. We denote for  $\mu \in \mathcal{M}(\mathcal{X})$  and  $\nu \in \mathcal{M}(\mathcal{Y})$ ,  $\mu \otimes \nu$  the tensor product of the measures  $\mu$  and  $\nu$ , and  $\mu \ll \nu$  means that  $\nu$  dominates  $\mu$ . We denote  $\Pi_1 : (x, y) \in \mathcal{X} \times \mathcal{Y} \mapsto x$  and

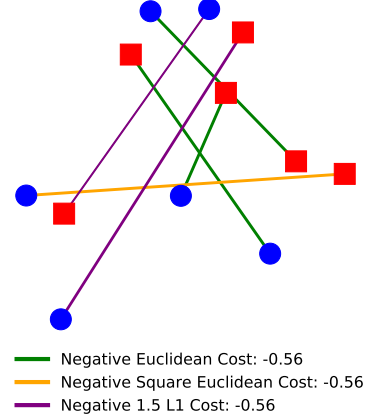


Figure 1: Equitable and optimal division of the resources between  $N = 3$  different negative costs (i.e. utilities) given by EOT. Utilities have been normalized. Blue dots and red squares represent the different elements of resources available in each cake. We consider the case where there is exactly one unit of supply per element in the cakes, which means that we consider uniform distributions. Note that the partition between the agents is equitable (i.e. utilities are equal) and proportional (i.e. utilities are larger than  $1/N$ ).

$\Pi_2 : (x, y) \in \mathcal{X} \times \mathcal{Y} \mapsto y$  respectively the projections on  $\mathcal{X}$  and  $\mathcal{Y}$ , which are continuous applications. For an application  $g$  and a measure  $\mu$ , we denote  $g_{\#}\mu$  the pushforward measure of  $\mu$  by  $g$ . For  $\mathcal{X}$  and  $\mathcal{Y}$  two Polish spaces, we denote  $\text{LSC}(\mathcal{X} \times \mathcal{Y})$  the space of lower semi-continuous functions on  $\mathcal{X} \times \mathcal{Y}$ ,  $\text{LSC}^+(\mathcal{X} \times \mathcal{Y})$  the space of non-negative lower semi-continuous functions on  $\mathcal{X} \times \mathcal{Y}$  and  $\text{LSC}_*^-(\mathcal{X} \times \mathcal{Y})$  the set of negative bounded below lower semi-continuous functions on  $\mathcal{X} \times \mathcal{Y}$ . We also denote  $\text{C}^+(\mathcal{X} \times \mathcal{Y})$  the space of non-negative continuous functions on  $\mathcal{X} \times \mathcal{Y}$  and  $\text{C}_*^-(\mathcal{X} \times \mathcal{Y})$  the set of negative continuous functions on  $\mathcal{X} \times \mathcal{Y}$ . Let  $N \geq 1$  be an integer and denote  $\Delta_N^+ := \{\lambda \in \mathbb{R}_+^N \text{ s.t. } \sum_{i=1}^N \lambda_i = 1\}$ , the probability simplex of  $\mathbb{R}^N$ . For two positive measures of same mass  $\mu \in \mathcal{M}_+(\mathcal{X})$  and  $\nu \in \mathcal{M}_+(\mathcal{Y})$ , we define the set of couplings with marginals  $\mu$  and  $\nu$ :

$$\Pi_{\mu, \nu} := \{\gamma \text{ s.t. } \Pi_{1\#}\gamma = \mu, \Pi_{2\#}\gamma = \nu\}.$$

We introduce the subset of  $(\mathcal{M}_+^1(\mathcal{X}) \times \mathcal{M}_+^1(\mathcal{Y}))^N$  representing marginal decomposition:

$$\Upsilon_{\mu, \nu}^N := \left\{ (\mu_i, \nu_i)_{i=1}^N \text{ s.t. } \sum_i \mu_i = \mu, \sum_i \nu_i = \nu \right. \\ \left. \text{and } \forall i, \mu_i(\mathcal{X}) = \nu_i(\mathcal{Y}) \right\}.$$

We also define the following subset of  $\mathcal{M}_+(\mathcal{X} \times \mathcal{Y})^N$  corresponding to the coupling decomposition:

$$\Gamma_{\mu, \nu}^N := \left\{ (\gamma_i)_{i=1}^N \text{ s.t. } \Pi_{1\#} \sum \gamma_i = \mu, \Pi_{2\#} \sum \gamma_i = \nu \right\}.$$

### 3.1 Primal Formulation

Consider a fair division problem where several agents aim to share two sets of resources,  $\mathcal{X}$  and  $\mathcal{Y}$ , and assume that there is a divisible amount of each resource  $x \in \mathcal{X}$  (resp.  $y \in \mathcal{Y}$ ) that is available. Formally, we consider the case where resources are no more sets but rather distributions on these sets. Denote  $\mu$  and  $\nu$  the distribution of resources on respectively  $\mathcal{X}$  and  $\mathcal{Y}$ . For example, one might think about a situation where agents want to share fruit juices and ice creams and there is a certain volume of each type of fruit juices and a certain mass of each type of ice creams available. Moreover each agent defines his or her paired preferences for each couple  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ . Formally, each person  $i$  is associated to an upper semi-continuous mapping  $u_i : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^+$  corresponding to his or her preference for any given pair  $(x, y)$ . For example, one may prefer to eat chocolate ice cream with apple juice, but may prefer pineapple juice when it comes with vanilla ice cream. The total utility for an individual  $i$  and a pairing  $\gamma_i \in \mathcal{M}_+(\mathcal{X} \times \mathcal{Y})$  is then given by  $V_i(\gamma_i) := \int u_i d\gamma_i$ . To partition fairly among individuals, we maximize the minimum of individual utilities.

From a transport point of view, let assume that there are  $N$  workers available to transport a distribution  $\mu$  to another one  $\nu$ . The cost of a worker  $i$  to transport a unit mass from location  $x$  to the location  $y$  is  $c_i(x, y)$ . To partition the work among the  $N$  workers fairly, we minimize the maximum of individual costs.

These problems are in fact the same where the utility  $u_i$ , defined in the fair division problem, might be interpreted as the opposite of the cost  $c_i$  defined in the transportation problem, i.e. for all  $i$ ,  $c_i = -u_i$ . The two above problem motivate the introduction of EOT defined as follows.

**Definition 1** (Equitable and Optimal Transport). *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be Polish spaces. Let  $\mathbf{c} := (c_i)_{1 \leq i \leq N}$  be a family of bounded below lower semi-continuous cost functions on  $\mathcal{X} \times \mathcal{Y}$ , and  $\mu \in \mathcal{M}_+^1(\mathcal{X})$  and  $\nu \in \mathcal{M}_+^1(\mathcal{Y})$ . We define the equitable and optimal transport primal problem:*

$$\text{EOT}_{\mathbf{c}}(\mu, \nu) := \inf_{(\gamma_i)_{i=1}^N \in \Gamma_{\mu, \nu}^N} \max_i \int c_i d\gamma_i. \quad (1)$$

We prove along with Theorem 1 that the problem is well defined and the infimum is attained. Lower-semi continuity is a standard assumption in OT. In fact, it is the weakest condition to prove Kantorovich duality (Villani, 2003, Chap. 1). Note that the problem defined here is a linear optimization problem and when  $N = 1$  we recover standard optimal transport. Figure 1 illustrates the equitable and optimal transport

problem we consider. Figure 5 in Appendix D shows an illustration with respect to the transport viewpoint in the exact same setting, i.e.  $c_i = -u_i$ . As expected, the couplings obtained in the two situations are not the same.

We now show that in fact, EOT optimum satisfies equality constraints in case of constant sign costs, i.e. total utility/cost of each individual are equal in the optimal partition. See Appendix A.2 for the proof.

**Proposition 1** (EOT solves the problem under equality constraints). *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be Polish spaces. Let  $\mathbf{c} := (c_i)_{1 \leq i \leq N} \in \text{LSC}^+(\mathcal{X} \times \mathcal{Y})^N \cup \text{LSC}_*^-(\mathcal{X} \times \mathcal{Y})^N$ ,  $\mu \in \mathcal{M}_+^1(\mathcal{X})$  and  $\nu \in \mathcal{M}_+^1(\mathcal{Y})$ . Then the following are equivalent:*

- $(\gamma_i^*)_{i=1}^N \in \Gamma_{\mu, \nu}^N$  is solution of Eq. (1),
- $(\gamma_i^*)_{i=1}^N \in \underset{(\gamma_i)_{i=1}^N \in \Gamma_{\mu, \nu}^N}{\text{argmin}} \left\{ t \text{ s.t. } \forall i \int c_i d\gamma_i = t \right\}.$

Moreover,

$$\text{EOT}_{\mathbf{c}}(\mu, \nu) = \min_{(\gamma_i)_{i=1}^N \in \Gamma_{\mu, \nu}^N} \left\{ t \text{ s.t. } \forall i \int c_i d\gamma_i = t \right\}.$$

This property highly relies on the sign of the costs. For instance if two costs are considered, one always positive and the other always negative, then the constraints cannot be satisfied. When the cost functions are non-negatives, EOT refers to a transportation problem while when the costs are all negatives, costs become utilities and EOT refers to a fair division problem. The two points of view are concordant, but proofs and interpretations rely on the sign of the costs.

### 3.2 An Equitable and Proportional Division

When the cost functions considered  $c_i$  are all negatives, EOT become a fair division problem where the utility functions are defined as  $u_i := -c_i$ . Indeed according to Proposition 1, EOT solves

$$\max_{(\gamma_i)_{i=1}^N \in \Gamma_{\mu, \nu}^N} \left\{ t \text{ s.t. } \forall i, \int u_i d\gamma_i = t \right\}.$$

Recall that in our model, the total utility of the agent  $i$  is given by  $V_i(\gamma_i) := \int u_i d\gamma_i$ . Therefore EOT aims to maximize the total utility of each agent  $i$  while ensuring that they are all equal. Let us now analyze which fairness conditions the partition induced by EOT verifies. Assume that the utilities are normalized, i.e.,  $\forall i$ , there exists  $\gamma_i \in \mathcal{M}_+^1(\mathcal{X} \times \mathcal{Y})$  such that  $V_i(\gamma_i) = 1$ . For example one might consider the cases where  $\forall i$ ,  $\gamma_i = \mu \otimes \nu$  or  $\gamma_i \in \underset{\gamma \in \Pi_{\mu, \nu}}{\text{argmin}} \int c_i d\gamma$ . Then any solution  $(\gamma_i^*)_{i=1}^N \in \Gamma_{\mu, \nu}^N$  of EOT satisfies:

- **Proportionality:** for all  $i$ ,  $V_i(\gamma_i^*) \geq 1/N$ ,

- **Equitability:** for all  $i, j$ ,  $V_i(\gamma_i^*) = V_j(\gamma_j^*)$ .

Proportionality is a standard fair division criterion for which a resource is divided among  $N$  agents, giving each agent at least  $1/N$  of the heterogeneous resource by his/her own subjective valuation. Therefore here, this situation corresponds to the case where the normalized utility of each agent is at least  $1/N$ . Moreover, an equitable division is a division of an heterogeneous resource, in which each partner is equally happy with his/her share. Here this corresponds to the case where the utility of each agent are all equal.

The problem solved by EOT is a fair division problem where heterogeneous resources have to be shared among multiple agents according to their preferences. This problem is a relaxation of the two cake-cutting problem when there are a divisible amount of each item of the cakes. In that case, cakes are distributions and EOT makes a proportional and equitable partition of them. Details are left in Appendix A.2.

**Fair Cake-cutting.** Consider the case where the cake is an heterogeneous resource and there is a certain divisible quantity of each type of resource available. For example chocolate and vanilla are two types of resource present in the cake for which a certain mass is available. In that case, each type of resource in the cake is pondered by the actual quantity present in the cake. Up to a normalization, the cake is no more the set  $\mathcal{X}$  but rather a distribution on this set. Note that for the two points of view to coincide, it suffices to assume that there is exactly the same amount of mass for each type of resources available in the cake. In that case, the cake can be represented by the uniform distribution over the set  $\mathcal{X}$ , or equivalently the set  $\mathcal{X}$  itself. When cakes are distributions, the fair cutting cake problem can be interpreted as a particular case of EOT when the utilities of the agents do not depend on the variable  $y \in \mathcal{Y}$ . In short, we consider that utilities are functions of the form  $u_i(x, y) = v_i(x)$  for all  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ . The normalization of utilities can be cast as follows:  $\forall i$ ,  $V_i(\mu) = \int v_i(x) d\mu(x) = 1$ . Then Proposition 1 shows that the partition of the cake made by EOT is proportional and equitable. Note that for EOT to coincide with the classical cake-cutting problem, one needs to consider that the uniform masses of the cake associated to each type of resource cannot be splitted. This can be interpreted as a Monge formulation (Villani, 2003) of EOT which is out of the scope of this paper.

### 3.3 Optimality of EOT

We next investigate the coupling obtained by solving EOT. In the next proposition, we show that under the same assumptions of Proposition 1, EOT solutions are

optimal transportation plans. See Appendix A.3 for the proof.

**Proposition 2** (EOT realizes optimal plans). *Under the same conditions of Proposition 1, for any  $(\gamma_i^*)_{i=1}^N \in \Gamma_{\mu, \nu}^N$  solution of Eq. (1), we have for all  $i \in \{1, \dots, N\}$*

$$\gamma_i^* \in \operatorname{argmin}_{\gamma \in \Pi_{\mu_i^*, \nu_i^*}} \int c_i d\gamma \quad (2)$$

$$\text{where } \mu_i^* := \Pi_{1\#} \gamma_i^*, \nu_i^* := \Pi_{2\#} \gamma_i^*,$$

and

$$\begin{aligned} \text{EOT}_{\mathbf{c}}(\mu, \nu) &= \min_{(\mu_i, \nu_i)_{i=1}^N \in \Upsilon_{\mu, \nu}^N} t \\ \text{s.t. } \forall i \quad W_{c_i}(\mu_i, \nu_i) &= t. \end{aligned} \quad (3)$$

Given the optimal matchings  $(\gamma_i^*)_{i=1}^N \in \Gamma_{\mu, \nu}^N$ , one can easily obtain the partition of the agents of each marginals. Indeed for all  $i$ ,  $\mu_i^* := \Pi_{1\#} \gamma_i^*$  and  $\nu_i^* := \Pi_{2\#} \gamma_i^*$  represent respectively the portion of the agent  $i$  from distributions  $\mu$  and  $\nu$ .

**Remark 1** (Utilitarian and Optimal Transport). *To contrast with EOT, an alternative problem is to maximize the sum of the total utilities of agents, or equivalently minimize the sum of the total costs of agents. This problem can be cast as follows:*

$$\inf_{(\gamma_i)_{i=1}^N \in \Gamma_{\mu, \nu}^N} \sum_i \int c_i d\gamma_i \quad (4)$$

Here one aims to maximize the total utility of all the agents, while in EOT we aim to maximize the total utility per agent under egalitarian constraint. The solution of (4) is not fair among agents and one can show that this problem is actually equal to  $W_{\min_i(c_i)}(\mu, \nu)$ . Details can be found in Appendix C.1.

### 3.4 Dual Formulation

Let us now introduce the dual formulation of the problem and show that strong duality holds under some mild assumptions. See Appendix A.4 for the proof.

**Theorem 1** (Strong Duality). *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be Polish spaces. Let  $\mathbf{c} := (c_i)_{i=1}^N$  be bounded below lower semi-continuous costs. Then strong duality holds, i.e. for  $(\mu, \nu) \in \mathcal{M}_+^1(\mathcal{X}) \times \mathcal{M}_+^1(\mathcal{Y})$ :*

$$\text{EOT}_{\mathbf{c}}(\mu, \nu) = \sup_{\substack{\lambda \in \Delta_N^+ \\ (f, g) \in \mathcal{F}_{\mathbf{c}}^\lambda}} \int f d\mu + \int g d\nu \quad (5)$$

where  $\mathcal{F}_{\mathbf{c}}^\lambda := \{(f, g) \in \mathcal{C}^b(\mathcal{X}) \times \mathcal{C}^b(\mathcal{Y}) \text{ s.t. } \forall i \in \{1, \dots, N\}, f \oplus g \leq \lambda_i c_i\}$ .

This theorem holds under the same hypothesis and follows the same reasoning as the one in (Villani, 2003,

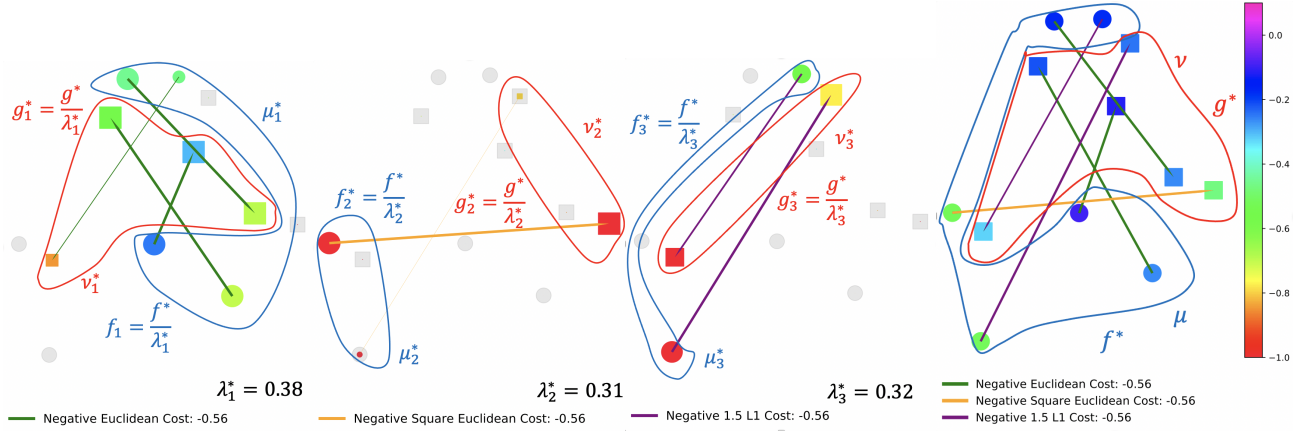


Figure 2: *Left, middle left, middle right*: the size of dots and squares is proportional to the weight of their representing atom in the distributions  $\mu_k^*$  and  $\nu_k^*$  respectively. The utilities  $f_k^*$  and  $g_k^*$  for each point in respectively  $\mu_k^*$  and  $\nu_k^*$  are represented by the color of dots and squares according to the color scale on the right hand side. The gray dots and squares correspond to the points that are ignored by agent  $k$  in the sense that there is no mass or almost no mass in distributions  $\mu_k^*$  or  $\nu_k^*$ . *Right*: the size of dots and squares are uniform since they correspond to the weights of uniform distributions  $\mu$  and  $\nu$  respectively. The values of  $f^*$  and  $g^*$  are given also by the color at each point. Note that each agent gets exactly the same total utility, corresponding exactly to EOT. This value can be computed using dual formulation (5) and for each figure it equals the sum of the values (encoded with colors) multiplied by the weight of each point (encoded with sizes).

Theorem 1.3). While the primal formulation of the problem is easy to understand, we want to analyse situations where the dual variables also play a role. For that purpose we show in the next proposition a simple characterisation of the primal-dual optimality in case of constant sign cost functions. See Appendix A.5 for the proof.

**Proposition 3.** *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be compact Polish spaces. Let  $\mathbf{c} := (c_i)_{1 \leq i \leq N} \in C^+(\mathcal{X} \times \mathcal{Y})^N \cup C^-(\mathcal{X} \times \mathcal{Y})^N$ ,  $\mu \in \mathcal{M}_+^1(\mathcal{X})$  and  $\nu \in \mathcal{M}_+^1(\mathcal{Y})$ . Let also  $(\gamma_k)_{k=1}^N \in \Gamma_{\mu, \nu}^N$  and  $(\lambda, f, g) \in \Delta_n^+ \times C^b(\mathcal{X}) \times C^b(\mathcal{Y})$ . Then Eq. (5) admits a solution and the following are equivalent:*

- $(\gamma_k)_{k=1}^N$  is a solution of Eq. (1) and  $(\lambda, f, g)$  is a solution of Eq. (5).
- 1.  $\forall i \in \{1, \dots, N\}, f \oplus g \leq \lambda_i c_i$
  2.  $\forall i, j \in \{1, \dots, N\} \int c_i d\gamma_i = \int c_j d\gamma_j$
  3.  $f \oplus g = \lambda_i c_i \quad \gamma_i$ -a.e.

**Remark 2.** *It is worth noting that when we assume that  $\mathbf{c} := (c_i)_{1 \leq i \leq N} \in C^+(\mathcal{X} \times \mathcal{Y})^N \cup C^-(\mathcal{X} \times \mathcal{Y})^N$ , then we can refine the second point of the equivalence presented in Proposition 3 by adding the following condition:  $\forall i \in \{1, \dots, N\} \lambda_i \neq 0$ .*

Given two distributions of resources represented by the measures  $\mu$  and  $\nu$ , and  $N$  utility functions denoted  $(u_i)_{i=1}^N$ , we want to find an *equitable* and *stable* partition among the agents in case of *transferable utilities*. Let  $k$  be an agent. We say that his or her utility is

transferable when once  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$  get matched, he or she has to decide how to split his or her associated utility  $u_k(x, y)$ . She or he divides  $u_k(x, y)$  into a quantity  $f_k(x)$  which can be seen as the utility of having  $x$  and  $g_k(y)$  for having  $y$ . Therefore in that problem we ask for  $(\gamma_k, f_k, g_k)_{k=1}^N$  such that

$$u_k(x, y) = f_k(x) + g_k(y) \quad \gamma_k\text{-a.e.} \quad (6)$$

Moreover, for the partition to be *stable* (Sotomayor and Roth, 1990), we want to ensure that, for every agent  $k$ , none of the resources  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$  that have not been matched together for this agent would increase their utilities,  $f_k(x)$  and  $g_k(y)$ , if there were matched together in the current matching instead. Formally we ask that for  $k \in \{1, \dots, N\}$  and all  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ ,

$$f_k(x) + g_k(y) \geq u_k(x, y). \quad (7)$$

Indeed if there exist  $k, x$  and  $y$  such that  $u_k(x, y) > f_k(x) + g_k(y)$ , then  $x$  and  $y$  will not be matched together in the share of the agent  $k$  and he can improve his utility for both  $x$  and  $y$  by matching  $x$  with  $y$ .

Finally we aim to share equitably the resources among the agents which boils down to ask

$$\forall i, j \in \{1, \dots, N\} \int u_i d\gamma_i = \int u_j d\gamma_j \quad (8)$$

Thanks to Proposition 3, finding  $(\gamma_k, f_k, g_k)_{k=1}^N$  satisfying (6), (7) and (8) can be done by solving Eq. (1)

and Eq. (5). Indeed let  $(\gamma_k)_{k=1}^N$  an optimal solution of Eq. (1) and  $(\lambda, f, g)$  an optimal solution of Eq. (5). Then by denoting for all  $k = 1, \dots, N$ ,  $f_k = \frac{f}{\lambda_k}$  and  $g_k = \frac{g}{\lambda_k}$ , we obtain that  $(\gamma_k, f_k, g_k)_{k=1}^N$  solves the *equitable* and *stable* partition problem in case of *transferable utilities*. Note that again, we end up with equality constraints for the optimal dual variables. Indeed, for all  $i, j \in \{1, \dots, N\}$ , at optimality we have  $\int f_i + g_i d\gamma_i = \int f_j + g_j d\gamma_j$ . Figure 2 illustrates this formulation of the problem with dual potentials. Figure 7 in Appendix D shows the dual solutions with respect to the transport viewpoint in the exact same setting, i.e.  $c_i = -u_i$ . Once again, the obtained solutions differ.

### 3.5 Link with other Probability Metrics

In this section, we provide some topological properties on the object defined by the EOT problem. In particular, we make links with other known probability metrics, such as Dudley and Wasserstein metrics and give a tight upper bound.

When  $N = 1$ , recall from the definition (1) that the problem considered is exactly the standard OT problem. Moreover any EOT problem with  $k \leq N$  costs can always be rewritten as a EOT problem with  $N$  costs. See Appendix C.2 for the proof. From this property, it is interesting to note that, for any  $N \geq 1$ , EOT generalizes standard Optimal Transport.

**Optimal Transport.** Given a cost function  $c$ , if we consider the problem EOT with  $N$  costs such that, for all  $i$ ,  $c_i = N \times c$  then, the problem  $\text{EOT}_{\mathbf{c}}$  is exactly  $\text{W}_c$ . See Appendix C.2 for the proof.

Now we have seen that all standard OT problems are sub-cases of the EOT problem, one may ask whether EOT can recover other families of metrics different from standard OT. Indeed we show that the EOT problem recovers an important family of IPMs with supremum taken over the space of  $\alpha$ -Hölder functions with  $\alpha \in (0, 1]$ . See Appendix A.6 for the proof.

**Proposition 4.** *Let  $\mathcal{X}$  be a Polish space. Let  $d$  be a metric on  $\mathcal{X}^2$  and  $\alpha \in (0, 1]$ . Denote  $c_1 = 2 \times \mathbf{1}_{x \neq y}$ ,  $c_2 = d^\alpha$  and  $\mathbf{c} := (c_1, (N-1) \times c_2, \dots, (N-1) \times c_2) \in \text{LSC}(\mathcal{X} \times \mathcal{X})^N$  then for any  $(\mu, \nu) \in \mathcal{M}_+^1(\mathcal{X}) \times \mathcal{M}_+^1(\mathcal{X})$*

$$\text{EOT}_{\mathbf{c}}(\mu, \nu) = \sup_{f \in B_{d^\alpha}(\mathcal{X})} \int_{\mathcal{X}} f d\mu - \int_{\mathcal{X}} f d\nu \quad (9)$$

where  $B_{d^\alpha}(\mathcal{X}) := \{f \in C^b(\mathcal{X}) : \|f\|_\infty + \|f\|_\alpha \leq 1\}$  and  $\|f\|_\alpha := \sup_{x \neq y} \frac{|f(x) - f(y)|}{d^\alpha(x, y)}$ .

**Dudley Metric.** When  $\alpha = 1$ , then for  $(\mu, \nu) \in \mathcal{M}_+^1(\mathcal{X}) \times \mathcal{M}_+^1(\mathcal{X})$ , we have

$$\text{EOT}_{\mathbf{c}}(\mu, \nu) = \text{EOT}_{(c_1, d)}(\mu, \nu) = \beta_d(\mu, \nu)$$

where  $\beta_d$  is the *Dudley Metric* (Dudley et al., 1966). In other words, the Dudley metric can be interpreted as an equitable and optimal transport between the measures with the trivial cost and a metric  $d$ . We acknowledge that Chizat et al. (2018) made a link between Unbalanced Optimal Transport and the “flat metric”, an IPM close to the Dudley metric, defined on the space  $\{f : \|f\|_\infty \leq 1, \|f\|_1 \leq 1\}$ .

**Weak Convergence.** When  $d$  is an unbounded metric on  $\mathcal{X}$ , it is well known that  $\text{W}_{d^p}$  with  $p \in (0, +\infty)$  metrizes a convergence a bit stronger than weak convergence (Villani, 2003, Chap. 7). A sufficient condition for Wasserstein distances to metrize weak convergence on the space of distributions is that the metric  $d$  is bounded. In contrast, metrics defined by Eq. (9) do not require such assumptions and  $\text{EOT}_{(\mathbf{1}_{x \neq y}, d^\alpha)}$  metrizes the weak convergence of probability measures (Villani, 2003, Chap. 1-7).

For an arbitrary choice of costs  $(c_i)_{1 \leq i \leq N}$ , we obtain a tight upper control of EOT and show how it is related to the OT problem associated to each cost involved. See Appendix A.7 for the proof.

**Proposition 5.** *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be Polish spaces. Let  $\mathbf{c} := (c_i)_{1 \leq i \leq N}$  be a family of nonnegative lower semi-continuous costs. For any  $(\mu, \nu) \in \mathcal{M}_+^1(\mathcal{X}) \times \mathcal{M}_+^1(\mathcal{Y})$*

$$\text{EOT}_{\mathbf{c}}(\mu, \nu) \leq \left( \sum_{i=1}^N \frac{1}{\text{W}_{c_i}(\mu, \nu)} \right)^{-1} \quad (10)$$

Proposition 5 means that the minimal cost to transport all goods under the constraint that all workers contribute equally is lower than the case where agents share equitably and optimally the transport with distributions  $\mu_i$  and  $\nu_i$  respectively proportional to  $\mu$  and  $\nu$ , which equals the harmonic sum written in Equation (10).

**Example.** *Applying the above result in the case of the Dudley metric recovers the following inequality (Sriperumbudur et al., 2012, Proposition 5.1)*

$$\beta_d(\mu, \nu) \leq \frac{\text{TV}(\mu, \nu) \text{W}_d(\mu, \nu)}{\text{TV}(\mu, \nu) + \text{W}_d(\mu, \nu)}.$$

## 4 Entropic Relaxation

In their original form, as proposed by Kantorovich (1942), Optimal Transport distances are not a natural fit for applied problems: they minimize a network flow problem, with a supercubic complexity ( $n^3 \log n$ ) (Tarjan, 1997). Following the work of Cuturi (2013), we propose an entropic relaxation of EOT, obtain its dual formulation and derive an efficient algorithm to compute an approximation of EOT.



#### 4.1 Primal-Dual Formulation

Let us first extend the notion of Kullback-Leibler divergence for positive Radon measures. Let  $\mathcal{Z}$  be a Polish space, for  $\mu, \nu \in \mathcal{M}_+(\mathcal{Z})$ , we define the generalized Kullback-Leibler divergence as  $\text{KL}(\mu||\nu) = \int \log \frac{d\mu}{d\nu} d\mu + \int d\nu - \int d\mu$  if  $\mu \ll \nu$ , and  $+\infty$  otherwise. We introduce the following regularized version of EOT.

**Definition 2** (Entropic relaxed primal problem). *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be two Polish spaces,  $\mathbf{c} := (c_i)_{1 \leq i \leq N}$  a family of bounded below lower semi-continuous costs lower semi-continuous costs on  $\mathcal{X} \times \mathcal{Y}$  and  $\boldsymbol{\varepsilon} := (\varepsilon_i)_{1 \leq i \leq N}$  be non negative real numbers. For  $(\mu, \nu) \in \mathcal{M}_+^1(\mathcal{X}) \times \mathcal{M}_+^1(\mathcal{Y})$ , we define the EOT regularized primal problem:*

$$\text{EOT}_{\mathbf{c}}^{\boldsymbol{\varepsilon}}(\mu, \nu) := \inf_{\gamma \in \Gamma_{\mu, \nu}^N} \max_i \int c_i d\gamma_i + \sum_{j=1}^N \varepsilon_j \text{KL}(\gamma_j || \mu \otimes \nu)$$

Note that here we sum the generalized Kullback-Leibler divergences since our objective is function of  $N$  measures in  $\mathcal{M}_+(\mathcal{X} \times \mathcal{Y})$ . This problem can be compared with the one from standard regularized OT. In the case where  $N = 1$ , we recover the standard regularized OT. For  $N \geq 1$ , the underlying problem is  $\sum_{i=1}^N \varepsilon_i$ -strongly convex. Moreover, we prove the essential property that as  $\boldsymbol{\varepsilon} \rightarrow 0$ , the regularized problem converges to the standard problem. See Appendix C.3 for the full statement and the proof. As a consequence, entropic regularization is a consistent approximation of the original problem we introduced in Section 3.1. Next theorem shows that strong duality holds for lower semi-continuous costs and compact spaces. This is the basis of the algorithm we will propose in Section 4.2. See Appendix A.8 for the proof.

**Theorem 2** (Duality for the regularized problem). *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be two compact Polish spaces,  $\mathbf{c} := (c_i)_{1 \leq i \leq N}$  a family of bounded below lower semi-continuous costs on  $\mathcal{X} \times \mathcal{Y}$  and  $\boldsymbol{\varepsilon} := (\varepsilon_i)_{1 \leq i \leq N}$  be non negative numbers. For  $(\mu, \nu) \in \mathcal{M}_+^1(\mathcal{X}) \times \mathcal{M}_+^1(\mathcal{Y})$ , strong duality holds:*

$$\text{EOT}_{\mathbf{c}}^{\boldsymbol{\varepsilon}}(\mu, \nu) = \sup_{\lambda \in \Delta_N^+} \sup_{\substack{f \in \mathcal{C}_b(\mathcal{X}) \\ g \in \mathcal{C}_b(\mathcal{Y})}} \int f d\mu + \int g d\nu - \sum_{i=1}^N \varepsilon_i \left( \int e^{\frac{f(x) + g(y) - \lambda_i c_i(x, y)}{\varepsilon_i}} d\mu(x) d\nu(y) - 1 \right) \quad (11)$$

and the infimum of the primal problem is attained.

As in standard regularized optimal transport there is a link between primal and dual variables at optimum. Let  $\gamma^*$  solving the regularized primal problem and  $(f^*, g^*, \lambda^*)$  solving the dual one:

$$\forall i, \gamma_i^* = \exp \left( \frac{f^* + g^* - \lambda_i^* c_i}{\varepsilon_i} \right) \cdot \mu \otimes \nu$$

#### 4.2 Proposed Algorithms

---

##### Algorithm 1 Projected Alternating Maximization

---

**Input:**  $\mathbf{C} = (C_i)_{1 \leq i \leq N}$ ,  $a, b, \varepsilon, L_\lambda$

**Init:**  $f^0 \leftarrow \mathbf{1}_n$ ;  $g^0 \leftarrow \mathbf{1}_m$ ;  $\lambda^0 \leftarrow (1/N, \dots, 1/N) \in \mathbb{R}^N$

**for**  $k = 1, 2, \dots$  **do**

$$\begin{aligned} K^k &\leftarrow \sum_{i=1}^N K_i^{\lambda^{k-1}}, \\ c_k &\leftarrow \langle f^{k-1}, K^k g^{k-1} \rangle, \quad f^k \leftarrow \frac{c_k a}{K^k g^{k-1}}, \\ d_k &\leftarrow \langle f^k, K^k g^{k-1} \rangle, \quad g^k \leftarrow \frac{d_k b}{(K^k)^T f^k}, \\ \lambda^k &\leftarrow \text{Proj}_{\Delta_N^+} \left( \lambda^{k-1} + \frac{1}{L_\lambda} \nabla_\lambda F_{\mathbf{C}}^\varepsilon(\lambda^{k-1}, f^k, g^k) \right). \end{aligned}$$

**end**

**Result:**  $\lambda, f, g$

---

We can now present algorithms obtained from entropic relaxation to approximately compute the solution of EOT. Let  $\mu = \sum_{i=1}^n a_i \delta_{x_i}$  and  $\nu = \sum_{j=1}^m b_j \delta_{y_j}$  be discrete probability measures where  $a \in \Delta_n^+$ ,  $b \in \Delta_m^+$ ,  $\{x_1, \dots, x_n\} \subset \mathcal{X}$  and  $\{y_1, \dots, y_m\} \subset \mathcal{Y}$ . Moreover for all  $i \in \{1, \dots, N\}$  and  $\lambda > 0$ , define  $\mathbf{C} := (C_i)_{1 \leq i \leq N} \in (\mathbb{R}^{n \times m})^N$  with  $C_i := (c_i(x_k, y_l))_{k, l}$  the  $N$  cost matrices and  $K_i^\lambda := \exp(-\lambda C_i / \varepsilon)$ . Assume that  $\varepsilon_1 = \dots = \varepsilon_N = \varepsilon$ . Compared to the standard regularized OT, the main difference here is that the problem contains an additional variable  $\lambda \in \Delta_N^+$ . When  $N = 1$ , one can use Sinkhorn algorithm. However when  $N \geq 2$ , we do not have a closed form for updating  $\lambda$  when the other variables of the problem are fixed. In order to enjoy from the strong convexity of the primal formulation, we consider instead the dual associated with the equivalent primal problem given when the additional trivial constraint  $\mathbf{1}_n^T (\sum_i P_i) \mathbf{1}_m = 1$  is considered. In that the dual obtained is

$$\begin{aligned} \widehat{\text{EOT}}_{\mathbf{C}}^{\boldsymbol{\varepsilon}}(a, b) &= \sup_{\substack{\lambda \in \Delta_N^+ \\ f \in \mathbb{R}^n, g \in \mathbb{R}^m}} \langle f, a \rangle + \langle g, b \rangle \\ &\quad - \varepsilon \left[ \log \left( \sum_i \langle e^{f/\varepsilon}, K_i^{\lambda_i} e^{g/\varepsilon} \rangle \right) + 1 \right] \end{aligned}$$

We show that the new objective obtained above is smooth w.r.t  $(\lambda, f, g)$ . See Appendix C.4 for the proof. One can apply the accelerated projected gradient ascent (Beck and Teboulle, 2009; Tseng, 2008) which enjoys an optimal convergence rate for first order methods of  $\mathcal{O}(k^{-2})$  for  $k$  iterations.

It is also possible to adapt Sinkhorn algorithm to our problem. See Algorithm 1. We denoted by  $\text{Proj}_{\Delta_N^+}$  the orthogonal projection on  $\Delta_N^+$  (Shalev-Shwartz and Singer, 2006), whose complexity is in  $\mathcal{O}(N \log N)$ . The smoothness constant in  $\lambda$  in the algorithm is  $L_\lambda = \max_i \|C_i\|_\infty^2 / \varepsilon$ . In practice Alg. 1 gives better results than the accelerated gradient descent. Note



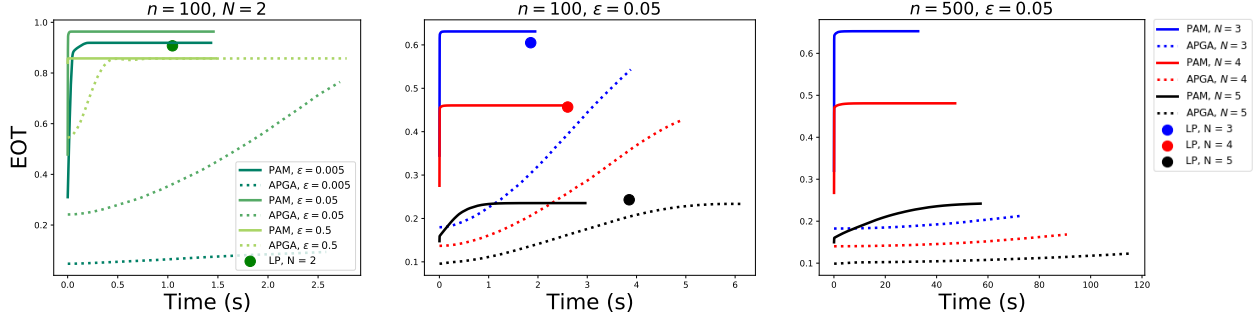


Figure 3: Comparison of the time-accuracy tradeoffs between the different proposed algorithms. *Left*: we consider the case where the number of days is  $N = 2$ , the size of support for both measures is  $n = m = 100$  and we vary  $\varepsilon$  from 0.005 to 0.5. *Middle*: we fix  $n = m = 100$  and the regularization  $\varepsilon = 0.05$  and we vary the number of days  $N$  from 3 to 5. *Right*: the setting considered is the same as in the figure in the middle, however we increase the sample size such that  $n = m = 500$ . Note that in that case, **LP** is too costly to be computed.

that the proposed algorithm differs from the Sinkhorn algorithm in many points and therefore the convergence rates cannot be applied here. Analyzing the rates of a *projected* alternating maximization method is, to the best of our knowledge, an unsolved problem. Further work will be devoted to study the convergence of this algorithm. We illustrate Algorithm 1 by showing the convergence of the regularized version of EOT towards the ground truth when  $\varepsilon \rightarrow 0$  in the case of the Dudley Metric. See Figure 8 in Appendix D.

## 5 Other applications of EOT

**Minimal Transportation Time.** Assume there are  $N$  internet service providers who propose different debits to transport data across locations, and one needs to transfer data from multiple servers to others, the fastest as possible. We assume that  $c_i(x, y) \geq 0$  corresponds to the transportation time needed by provider  $i$  to transport one unit of data from a server  $x$  to a server  $y$ . For instance, the unit of data can be one Megabit. Then  $\int c_i d\gamma_i$  corresponds to the time taken by provider  $i$  to transport  $\mu_i = \Pi_{1\#}\gamma_i$  to  $\nu_i = \Pi_{2\#}\gamma_i$ . Assuming the transportation can be made in parallel and given a partition of the transportation task  $(\gamma_i)_{i=1}^N$ ,  $\max_i \int c_i d\gamma_i$  corresponds to the total time of transport the data  $\mu = \Pi_{1\#} \sum \gamma_i$  to the locations  $\nu = \Pi_{2\#} \sum \gamma_i$  according to this partition. Then EOT, which minimizes  $\max_i \int c_i d\gamma_i$ , is finding the fastest way to transport the data from  $\mu$  to  $\nu$  by splitting the task among the  $N$  internet service providers. Note that at optimality, all the internet service providers finish their transportation task at the same time (see Proposition 1).

**Sequential Optimal Transport.** Consider the situation where an agent aims to transport goods from some stocks to some stores in the next  $N$  days. The cost to transport one unit of good from a stock located

at  $x$  to a store located at  $y$  may vary across the days. For example the cost of transportation may depend on the price of gas, or the daily weather conditions. Assuming that he or she has a good knowledge of the daily costs of the  $N$  coming days, he or she may want a transportation strategy such that his or her daily cost is as low as possible. By denoting  $c_i$  the cost of transportation the  $i$ -th day, and given a strategy  $(\gamma_i)_{i=1}^N$ , the maximum daily cost is then  $\max_i \int c_i d\gamma_i$ , and EOT therefore finds the cheapest strategy to spread the transport task in the next  $N$  days such that the maximum daily cost is minimized. Note that at optimality he or she has to spend the exact same amount everyday.

In Figure 3 we aim to simulate the Sequential OT problem and compare the time-accuracy trade-offs of the proposed algorithms. Let us consider a situation where one wants to transport merchandises from  $\mu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$  to  $\nu = \frac{1}{m} \sum_{j=1}^m \delta_{y_j}$  in  $N$  days. Here we model the locations  $\{x_i\}$  and  $\{y_j\}$  by drawing them independently from two Gaussian distributions in  $\mathbb{R}^2$ :  $\forall i, x_i \sim \mathcal{N}(\begin{pmatrix} 3 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix})$  and  $\forall j, y_j \sim \mathcal{N}(\begin{pmatrix} 4 \\ -1.2 \end{pmatrix}, \begin{pmatrix} 1 & -2 \\ -2 & 1 \end{pmatrix})$ . We assume that everyday there is wind modeled by a vector  $w \sim \mathcal{U}(B(0, 1))$  where  $B(0, 1)$  is the unit ball in  $\mathbb{R}^2$  that is perfectly known in advance. We define the cost of transportation on day  $i$  as  $c_i(x, y) = \|y - x\| - 0.7\langle w_i, y - x \rangle$  to model the effect of the wind on the transportation cost. In the following figures we plot the estimates of EOT obtained from the proposed algorithms in function of the runtime for various sample sizes  $n$ , number of days  $N$  and regularizations  $\varepsilon$ . **PAM** denotes Alg. 1, **APGA** denotes Alg. 2 (See Appendix C.4), **LP** denotes the linear program which solves exactly the primal formulation of the EOT problem. Note that when **LP** is computable (i.e.  $n \leq 100$ ), it is therefore the ground truth. We show that in all the settings, **PAM** performs better than **APGA** and provides very high accuracy with order of magnitude faster than **LP**.

## Acknowledgments

The authors would like to thank V. Do and Y. Chevalerey for fruitful discussions. We gratefully acknowledge support from "Chaire d'excellence de l'IDEX Paris Saclay".

## References

- David Alvarez-Melis, Stefanie Jegelka, and Tommi S Jaakkola. Towards optimal transport with global invariances. *arXiv preprint 1806.09277*, 2018.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.
- Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1): 183–202, 2009.
- Dimitris Bertsimas and John N Tsitsiklis. *Introduction to Linear Optimization*. Athena Scientific, 1997.
- Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- Felix Brandt, Vincent Conitzer, Ulle Endriss, Jérôme Lang, and Ariel D Procaccia. *Handbook of computational social choice*. Cambridge University Press, 2016.
- Haim Brezis. *Functional analysis, Sobolev spaces and partial differential equations*. Springer Science & Business Media, 2010.
- Lénaïc Chizat, Gabriel Peyré, Bernhard Schmitzer, and François-Xavier Vialard. Unbalanced optimal transport: Dynamic and kantorovich formulations. *Journal of Functional Analysis*, 274(11):3090–3123, 2018.
- John Cloutier, Kathryn L Nyman, and Francis Edward Su. Two-player envy-free multi-cake division. *Mathematical Social Sciences*, 59(1):26–37, 2010.
- Nicolas Courty, Rémi Flamary, Devis Tuia, and Alain Rakotomamonjy. Optimal transport for domain adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 39(9):1853–1865, 2016.
- Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in neural information processing systems*, pages 2292–2300, 2013.
- Lester E Dubins and Edwin H Spanier. How to cut a cake fairly. *The American Mathematical Monthly*, 68 (1P1):1–17, 1961.
- Richard M. Dudley. The speed of mean Glivenko-Cantelli convergence. *Annals of Mathematical Statistics*, 40(1):40–50, 1969.
- Richard Mansfield Dudley et al. Weak convergence of probabilities on nonseparable metric spaces and empirical measures on euclidean spaces. *Illinois Journal of Mathematics*, 10(1):109–126, 1966.
- Paul Dupuis and Richard S Ellis. *A weak convergence approach to the theory of large deviations*, volume 902. John Wiley & Sons, 2011.
- Arnaud Dupuy, Alfred Galichon, and Yifei Sun. Estimating matching affinity matrix under low-rank constraints. *arXiv preprint arXiv:1612.09585*, 2016.
- Gintare Karolina Dziugaite, Daniel M Roy, and Zoubin Ghahramani. Training generative neural networks via maximum mean discrepancy optimization. *arXiv preprint arXiv:1505.03906*, 2015.
- Jean Feydy, Thibault Séjourné, François-Xavier Vialard, Shun-Ichi Amari, Alain Trounev, and Gabriel Peyré. Interpolating between optimal transport and mmd using sinkhorn divergences. *arXiv preprint arXiv:1810.08278*, 2018.
- Nicolas Fournier and Arnaud Guillin. On the rate of convergence in Wasserstein distance of the empirical measure. *Probability Theory and Related Fields*, 162 (3-4):707–738, 2015.
- Aude Genevay, Gabriel Peyré, and Marco Cuturi. Learning generative models with sinkhorn divergences, 2017.
- Aude Genevay, Lénaïc Chizat, Francis Bach, Marco Cuturi, and Gabriel Peyré. Sample complexity of sinkhorn divergences. *arXiv preprint arXiv:1810.02733*, 2018.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar):723–773, 2012.
- Hisham Husain, Richard Nock, and Robert C Williamson. A primal-dual link between gans and autoencoders. In *Advances in Neural Information Processing Systems*, pages 413–422, 2019.
- Hicham Janati, Thomas Bazeille, Bertrand Thirion, Marco Cuturi, and Alexandre Gramfort. Multi-subject meg/eeg source imaging with sparse multi-task regression. *NeuroImage*, page 116847, 2020.
- Leonid Kantorovich. On the transfer of masses (in russian). *Doklady Akademii Nauk*, 37(2):227–229, 1942.
- Tor Lattimore, Koby Crammer, and Csaba Szepesvári. Linear multi-resource allocation with semi-bandit

- feedback. In *Advances in Neural Information Processing Systems*, pages 964–972, 2015.
- Ruilin Li, Xiaojing Ye, Haomin Zhou, and Hongyuan Zha. Learning to match via inverse optimal transport. *J. Mach. Learn. Res.*, 20:80–1, 2019.
- Erika Mackin and Lirong Xia. Allocating indivisible items in categorized domains. *arXiv preprint arXiv:1504.05932*, 2015.
- Hervé Moulin. *Fair division and collective welfare*. MIT press, 2004.
- Youssef Mroueh and Tom Sercu. Fisher gan. In *Advances in Neural Information Processing Systems*, pages 2513–2523, 2017.
- Alfred Müller. Integral probability metrics and their generating classes of functions. *Advances in Applied Probability*, 29(2):429–443, 1997.
- Yu Nesterov. Smooth minimization of non-smooth functions. *Mathematical programming*, 103(1):127–152, 2005.
- François-Pierre Paty and Marco Cuturi. Subspace robust wasserstein distances. *arXiv preprint arXiv:1901.08949*, 2019.
- Mathis Petrovich, Chao Liang, Yanbin Liu, Yao-Hung Hubert Tsai, Linchao Zhu, Yi Yang, Ruslan Salakhutdinov, and Makoto Yamada. Feature robust optimal transport for high-dimensional data. *arXiv preprint arXiv:2005.12123*, 2020.
- Tim Salimans, Han Zhang, Alec Radford, and Dimitris Metaxas. Improving GANs using optimal transport. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rkQkBNJAb>.
- M. Scetbon and G. Varoquaux. Comparing distributions:  $\ell_1$  geometry improves kernel two-sample testing, 2019.
- Meyer Scetbon and Marco Cuturi. Linear time sinkhorn divergences using positive features, 2020.
- 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.
- Shai Shalev-Shwartz and Yoram Singer. Efficient learning of label ranking by soft projections onto polyhedra. *Journal of Machine Learning Research*, 7(Jul):1567–1599, 2006.
- Maurice Sion. On general minimax theorems. *Pacific J. Math.*, 8(1):171–176, 1958. URL <https://projecteuclid.org:443/euclid.pjm/1103040253>.
- Marilda Sotomayor and Alvin Roth. Two-sided matching: A study in game-theoretic modelling and analysis. *Econometric Society Monographs*, (18), 1990.
- Bharath K Sriperumbudur, Kenji Fukumizu, Arthur Gretton, Bernhard Schölkopf, Gert RG Lanckriet, et al. On the empirical estimation of integral probability metrics. *Electronic Journal of Statistics*, 6:1550–1599, 2012.
- Ton Steerneman. On the total variation and hellinger distance between signed measures; an application to product measures. *Proceedings of the American Mathematical Society*, 88(4):684–688, 1983.
- H. Steinhaus. Sur la division pragmatique. *Econometrica*, 17:315–319, 1949. ISSN 00129682, 14680262. URL <http://www.jstor.org/stable/1907319>.
- Haodong Sun, Haomin Zhou, Hongyuan Zha, and Xiaojing Ye. Learning cost functions for optimal transport. *arXiv preprint arXiv:2002.09650*, 2020.
- Robert E. Tarjan. Dynamic trees as search trees via euler tours, applied to the network simplex algorithm. *Mathematical Programming*, 78(2):169–177, 1997.
- Paul Tseng. On accelerated proximal gradient methods for convex-concave optimization. *submitted to SIAM Journal on Optimization*, 1, 2008.
- Cédric Villani. *Topics in optimal transportation*. Number 58. American Mathematical Soc., 2003.
- Haibin Wang, Sujoy Sikdar, Xiaoxi Guo, Lirong Xia, Yongzhi Cao, and Hanpin Wang. Multi-type resource allocation with partial preferences. *arXiv preprint arXiv:1906.06836*, 2019.
- Weiran Wang and Miguel A. Carreira-Perpinan. Projection onto the probability simplex: An efficient algorithm with a simple proof, and an application, 2013.
- Eric Wong, Frank R. Schmidt, and J. Zico Kolter. Wasserstein adversarial examples via projected sinkhorn iterations, 2019.
- Karren Dai Yang, Karthik Damodaran, Saradha Venkatachalapathy, Ali C Soylemezoglu, GV Shivashankar, and Caroline Uhler. Predicting cell lineages using autoencoders and optimal transport. *PLoS computational biology*, 16(4):e1007828, 2020.

## Supplementary material

### A Proofs

#### A.1 Notations

Let  $\mathcal{Z}$  be a Polish space, we denote  $\mathcal{M}(\mathcal{Z})$  the set of Radon measures on  $\mathcal{Z}$  endowed with total variation norm:  $\|\mu\|_{TV} = \mu_+(\mathcal{Z}) + \mu_-(\mathcal{Z})$  with  $(\mu_+, \mu_-)$  is the Dunford decomposition of the signed measure  $\mu$ . We call  $\mathcal{M}_+(\mathcal{Z})$  the sets of positive Radon measures, and  $\mathcal{M}_+^1(\mathcal{Z})$  the set of probability measures. We denote  $\mathcal{C}^b(\mathcal{Z})$  the vector space of bounded continuous functions on  $\mathcal{Z}$  endowed with  $\|\cdot\|_\infty$  norm. We recall the *Riesz-Markov theorem*: if  $\mathcal{Z}$  is compact,  $\mathcal{M}(\mathcal{Z})$  is the topological dual of  $\mathcal{C}^b(\mathcal{Z})$ . Let  $\mathcal{X}$  and  $\mathcal{Y}$  be two Polish spaces. It is immediate that  $\mathcal{X} \times \mathcal{Y}$  is a Polish space. We denote for  $\mu \in \mathcal{M}(\mathcal{X})$  and  $\nu \in \mathcal{M}(\mathcal{Y})$ ,  $\mu \otimes \nu$  the tensor product of the measures  $\mu$  and  $\nu$ , and  $\mu \ll \nu$  means that  $\nu$  dominates  $\mu$ . We denote  $\Pi_1 : (x, y) \in \mathcal{X} \times \mathcal{Y} \mapsto x$  and  $\Pi_2 : (x, y) \in \mathcal{X} \times \mathcal{Y} \mapsto y$  respectively the projections on  $\mathcal{X}$  and  $\mathcal{Y}$ , which are continuous applications. For an application  $g$  and a measure  $\mu$ , we denote  $g_\# \mu$  the pushforward measure of  $\mu$  by  $g$ . For  $f : \mathcal{X} \rightarrow \mathbb{R}$  and  $g : \mathcal{Y} \rightarrow \mathbb{R}$ , we denote  $f \oplus g : (x, y) \in \mathcal{X} \times \mathcal{Y} \mapsto f(x) + g(y)$  the tensor sum of  $f$  and  $g$ . For  $\mathcal{X}$  and  $\mathcal{Y}$  two Polish spaces, we denote  $\text{LSC}(\mathcal{X} \times \mathcal{Y})$  the space of lower semi-continuous functions on  $\mathcal{X} \times \mathcal{Y}$ ,  $\text{LSC}^+(\mathcal{X} \times \mathcal{Y})$  the space of non-negative lower semi-continuous functions on  $\mathcal{X} \times \mathcal{Y}$  and  $\text{LSC}_*^-(\mathcal{X} \times \mathcal{Y})$  the set of negative bounded below lower semi-continuous functions on  $\mathcal{X} \times \mathcal{Y}$ . Let  $N \geq 1$  be an integer and denote  $\Delta_N^+ := \{\lambda \in \mathbb{R}_+^N \text{ s.t. } \sum_{i=1}^N \lambda_i = 1\}$ , the probability simplex of  $\mathbb{R}^N$ . For two positive measures of same mass  $\mu \in \mathcal{M}_+(\mathcal{X})$  and  $\nu \in \mathcal{M}_+(\mathcal{Y})$ , we define the set of couplings with marginals  $\mu$  and  $\nu$ :

$$\Pi_{\mu, \nu} := \{\gamma \text{ s.t. } \Pi_{1\#} \gamma = \mu, \Pi_{2\#} \gamma = \nu\}.$$

For  $\mu \in \mathcal{M}_+^1(\mathcal{X})$  and  $\nu \in \mathcal{M}_+^1(\mathcal{Y})$ , we introduce the subset of  $(\mathcal{M}_+^1(\mathcal{X}) \times \mathcal{M}_+^1(\mathcal{Y}))^N$  representing marginal decomposition:

$$\Upsilon_{\mu, \nu}^N := \{(\mu_i, \nu_i)_{i=1}^N \text{ s.t. } \sum_i \mu_i = \mu, \sum_i \nu_i = \nu \text{ and } \forall i, \mu_i(\mathcal{X}) = \nu_i(\mathcal{Y})\}.$$

We also define the following subset of  $\mathcal{M}_+(\mathcal{X} \times \mathcal{Y})^N$  corresponding to the coupling decomposition:

$$\Gamma_{\mu, \nu}^N := \left\{ (\gamma_i)_{i=1}^N \text{ s.t. } \Pi_{1\#} \sum_i \gamma_i = \mu, \Pi_{2\#} \sum_i \gamma_i = \nu \right\}.$$

#### A.2 Proof of Proposition 1

**Proof.** First, it is clear that  $\text{EOT}_{\mathbf{c}}(\mu, \nu) \geq \inf_{\gamma \in \Gamma_{\mu, \nu}^N} \{t \text{ s.t. } \forall i, t = \int c_i d\gamma_i\}$ . Let us now show that in fact it is an equality. Thanks to Theorem 1, the infimum is attained for  $\inf_{\gamma \in \Gamma_{\mu, \nu}^N} \max_i \int c_i d\gamma_i$ . Indeed recall that  $\Gamma_{\mu, \nu}^N$  is compact and that the objective is lower semi-continuous. Let  $\gamma^*$  be such a minimizer. Let  $I$  be the set of indices  $i$  such that  $\int c_i d\gamma_i^* = \text{EOT}_{\mathbf{c}}(\mu, \nu)$ . Assume that there exists  $j$  such that,  $\text{EOT}_{\mathbf{c}}(\mu, \nu) > \int c_j d\gamma_j^*$ .

In case of costs of  $\text{LSC}^+(\mathcal{X} \times \mathcal{Y})$ , for all  $i \in I$ , there exists  $(x_i, y_i) \in \text{Supp}(\gamma_i^*)$  such that  $c_i(x_i, y_i) > 0$ . Let us denote  $A_{(x_i, y_i)}$  measurable sets such that  $(x_i, y_i) \in A_{(x_i, y_i)}$  and let us denote  $\tilde{\gamma}$  defined as for all  $k \notin I \cup \{j\}$ ,  $\tilde{\gamma}_k = \gamma_k^*$ , for  $i \in I$ ,  $\tilde{\gamma}_i = \gamma_i^* - \epsilon \mathbf{1}_{A_{(x_i, y_i)}} \gamma_i^*$  and  $\tilde{\gamma}_j = \gamma_j^* + \sum_{i \in I} \epsilon \mathbf{1}_{A_{(x_i, y_i)}} \gamma_i^*$  for  $\epsilon$  sufficiently small so that  $\tilde{\gamma} \in \Gamma_{\mu, \nu}^N$ . Now,  $\max_k \int c_k d\tilde{\gamma}_k^* > \max_k \int c_k d\gamma_k^*$ , which contradicts that  $\gamma^*$  is a minimizer. Then for  $i, j$ ,  $\int c_i d\gamma_i^* = \int c_j d\gamma_j^*$ . And then:  $\text{EOT}_{\mathbf{c}}(\mu, \nu) = \inf_{\gamma \in \Gamma_{\mu, \nu}^N} \max_i \int c_i d\gamma_i$ .

In case of costs in  $\text{LSC}_*^-(\mathcal{X} \times \mathcal{Y})$ , there exists  $(x_0, y_0) \in \text{Supp}(\gamma_j^*)$  such that  $c_j(x_0, y_0) < 0$ . Let us denote  $A_{(x_0, y_0)}$  a measurable set such that  $(x_0, y_0) \in A_{(x_0, y_0)}$  and let us denote  $\tilde{\gamma}$  defined as for all  $k \notin I \cup \{j\}$ ,  $\tilde{\gamma}_k = \gamma_k^*$  and for all  $i \in I$ ,  $\tilde{\gamma}_i = \gamma_i^* + \frac{\epsilon}{|I|} \mathbf{1}_{A_{(x_0, y_0)}} \gamma_j^*$  and  $\tilde{\gamma}_j = \gamma_j^* - \epsilon \mathbf{1}_{A_{(x_0, y_0)}} \gamma_j^*$  for  $\epsilon$  sufficiently small so that  $\tilde{\gamma} \in \Gamma_{\mu, \nu}^N$ . Now,  $\max_k \int c_k d\tilde{\gamma}_k^* > \max_k \int c_k d\gamma_k^*$ , which contradicts that  $\gamma^*$  is a minimizer. Then for  $i, j$ ,  $\int c_i d\gamma_i^* = \int c_j d\gamma_j^*$ . And then:  $\text{EOT}_{\mathbf{c}}(\mu, \nu) = \inf_{\gamma \in \Gamma_{\mu, \nu}^N} \max_i \int c_i d\gamma_i$ .

It is clear that equitability is verified thanks to the previous proof. For proportionality, assume the normalization:  $\forall i$ , there exists  $\gamma_i \in \mathcal{M}_+^1(\mathcal{X} \times \mathcal{Y})$  such that  $V_i(\gamma_i) = 1$ . Then for each  $i$ ,  $V_i(\gamma_i/N) = 1/N$  and  $(\gamma_i)_i \in \Gamma_{\mu, \nu}^N$ . Then at optimum:  $\forall i$ ,  $V_i(\gamma_i^*) \geq 1/N$  and proportionality is verified.

### A.3 Proof of Proposition 2

**Proof.** We prove along with Theorem 1 that the infimum defining  $\text{EOT}_c(\mu, \nu)$  is attained. Let  $\gamma^*$  be this infimum. Then at optimum we have shown that for all  $i, j$ ,  $\int c_i d\gamma_i^* = \int c_j d\gamma_j^* = t$ . Let denote for all  $i$ ,  $\mu_i = \Pi_{1\#} \gamma_i^*$  and  $\nu_i = \Pi_{2\#} \gamma_i^*$ .

Let assume there exists  $i$  such that  $\int c_i d\gamma_i^* > W_{c_i}(\mu_i, \nu_i)$ . Let  $\gamma'_i$  realising the infimum of  $W_{c_i}(\mu_i, \nu_i)$ . Let  $\epsilon > 0$  be sufficiently small, then let define  $\tilde{\gamma}$  as follows: for all  $j \neq i$ ,  $\tilde{\gamma}_j = (1 - \epsilon)\gamma_j^*$ . and  $\tilde{\gamma}_i = \gamma'_i + \epsilon \sum_{j \neq i} \gamma_j^*$ . Then for all  $j \neq i$ ,  $\int c_j d\tilde{\gamma}_j = (1 - \epsilon)t$  and  $\int c_i d\tilde{\gamma}_i = W_{c_i}(\mu_i, \nu_i) + \epsilon \sum_{j \neq i} \int c_i d\gamma_j^*$ . It is clear that  $\tilde{\gamma} \in \Gamma_{\mu, \nu}^N$ . For  $\epsilon > 0$  sufficiently small,  $\max_i \int c_i d\tilde{\gamma}_i = (1 - \epsilon)t < t$ , which contradicts the optimality of  $\gamma^*$ .

A possible reformulation for EOT is:

$$\text{EOT}_c(\mu, \nu) = \min_{\substack{(\mu_i, \nu_i)_{i=1}^N \in \Upsilon_{\mu, \nu}^N \\ \forall i, \gamma_i \in \Pi_{\mu, \nu}}} \left\{ t \text{ s.t. } \int c_i d\gamma_i = t \right\}$$

We previously show that at optimum the couplings are optimal transport plans, then:

$$\text{EOT}_c(\mu, \nu) = \min_{(\mu_i, \nu_i)_{i=1}^N \in \Upsilon_{\mu, \nu}^N} \{ t \text{ s.t. } \forall i, W_{c_i}(\mu_i, \nu_i) = t \}$$

which concludes the proof.

### A.4 Proof of Theorem 1

To prove this theorem, one need to prove the three following technical lemmas. The first one shows the weak compacity of  $\Gamma_{\mu, \nu}^N$ .

**Lemma 1.** Let  $\mathcal{X}$  and  $\mathcal{Y}$  be Polish spaces, and  $\mu$  and  $\nu$  two probability measures respectively on  $\mathcal{X}$  and  $\mathcal{Y}$ . Then  $\Gamma_{\mu, \nu}^N$  is sequentially compact for the weak topology induced by  $\|\gamma\| = \max_{i=1, \dots, N} \|\gamma_i\|_{\text{TV}}$ .

**Proof.** Let  $(\gamma^n)_{n \geq 0}$  a sequence in  $\Gamma_{\mu, \nu}^N$ , and let us denote for all  $n \geq 0$ ,  $\gamma^n = (\gamma_i^n)_{i=1}^N$ . We first remark that for all  $i \in \{1, \dots, N\}$  and  $n \geq 0$ ,  $\|\gamma_i^n\|_{\text{TV}} \leq 1$  therefore for all  $i \in \{1, \dots, N\}$ ,  $(\gamma_i^n)_{n \geq 0}$  is uniformly bounded. Moreover as  $\{\mu\}$  and  $\{\nu\}$  are tight, for any  $\delta > 0$ , there exist  $K \subset \mathcal{X}$  and  $L \subset \mathcal{Y}$  compact sets such that

$$\mu(K^c) \leq \frac{\delta}{2} \quad \text{and} \quad \nu(L^c) \leq \frac{\delta}{2}. \quad (12)$$

Therefore, we obtain that for any for all  $i \in \{1, \dots, N\}$ ,

$$\gamma_i^n(K^c \times L^c) \leq \sum_{k=1}^N \gamma_k^n(K^c \times L^c) \quad (13)$$

$$\leq \sum_{k=1}^N \gamma_k^n(K^c \times \mathcal{Y}) + \gamma_k^n(\mathcal{X} \times L^c) \quad (14)$$

$$\leq \mu(K^c) + \nu(L^c) = \delta. \quad (15)$$

Therefore, for all  $i \in \{1, \dots, N\}$ ,  $(\gamma_i^n)_{n \geq 0}$  is tight and uniformly bounded and Prokhorov's theorem (Dupuis and Ellis, 2011, Theorem A.3.15) guarantees for all  $i \in \{1, \dots, N\}$ ,  $(\gamma_i^n)_{n \geq 0}$  admits a weakly convergent subsequence. By extracting a common convergent subsequence, we obtain that  $(\gamma^n)_{n \geq 0}$  admits a weakly convergent subsequence. By continuity of the projection, the limit also lives in  $\Gamma_{\mu, \nu}^N$  and the result follows.

Next lemma generalizes Rockafellar-Fenchel duality to our case.

**Lemma 2.** Let  $V$  be a normed vector space and  $V^*$  its topological dual. Let  $V_1, \dots, V_N$  be convex functions and lower semi-continuous on  $V$  and  $E$  a convex function on  $V$ . Let  $V_1^*, \dots, V_N^*, E^*$  be the Fenchel-Legendre transforms of  $V_1, \dots, V_N, E$ . Assume there exists  $z_0 \in V$  such that for all  $i$ ,  $V_i(z_0) < \infty$ ,  $E(z_0) < \infty$ , and for all  $i$ ,  $V_i$  is continuous at  $z_0$ . Then:

$$\inf_{u \in V} \sum_i V_i(u) + E(u) = \sup_{\substack{\gamma_1, \dots, \gamma_N, \gamma \in V^* \\ \sum_i \gamma_i = \gamma}} - \sum_i V_i^*(-\gamma_i) - E^*(\gamma)$$

**Proof.** This Lemma is an immediate application of Rockafellar-Fenchel duality theorem (Brezis, 2010, Theorem 1.12) and of Fenchel-Moreau theorem (Brezis, 2010, Theorem 1.11). Indeed,  $V = \sum_{i=1}^N V_i(u)$  is a convex function, lower semi-continuous and its Legendre-Fenchel transform is given by:

$$V^*(\gamma^*) = \inf_{\sum_{i=1}^N \gamma_i^* = \gamma^*} \sum_{i=1}^N V_i^*(\gamma_i^*). \quad (16)$$

Last lemma is an application of Sion's Theorem to this problem.

**Lemma 3.** Let  $\mathcal{X}$  and  $\mathcal{Y}$  be Polish spaces. Let  $\mathbf{c} = (c_i)_{1 \leq i \leq N}$  be a family of bounded below lower semi-continuous costs on  $\mathcal{X} \times \mathcal{Y}$ , then for  $\mu \in \mathcal{M}_+^1(\mathcal{X})$  and  $\nu \in \mathcal{M}_+^1(\mathcal{Y})$ , we have

$$\text{EOT}_{\mathbf{c}}(\mu, \nu) = \sup_{\lambda \in \Delta_N^+} \inf_{\gamma \in \Gamma_{\mu, \nu}^N} \sum_{i=1}^N \lambda_i \int_{\mathcal{X} \times \mathcal{Y}} c_i(x, y) d\gamma_i(x, y) \quad (17)$$

and the infimum is attained.

**Proof.** Taking for granted that a minmax principle can be invoked, we have

$$\begin{aligned} \sup_{\lambda \in \Delta_N^+} \inf_{\gamma \in \Gamma_{\mu, \nu}^N} \sum_{i=1}^N \lambda_i \int_{\mathcal{X} \times \mathcal{Y}} c_i(x, y) d\gamma_i(x, y) &= \inf_{\gamma \in \Gamma_{\mu, \nu}^N} \sup_{\lambda \in \Delta_N^+} \sum_{i=1}^N \lambda_i \int_{\mathcal{X} \times \mathcal{Y}} c_i(x, y) d\gamma_i(x, y) \\ &= \text{EOT}_{\mathbf{c}}(\mu, \nu) \end{aligned}$$

But thanks to Lemma 1, we have that  $\Gamma_{\mu, \nu}^N$  is compact for the weak topology. And  $\Delta_N^+$  is convex. Moreover the objective function  $f : (\lambda, \gamma) \in \Delta_N^+ \times \Gamma_{\mu, \nu}^N \mapsto \sum_{i=1}^N \lambda_i \int_{\mathcal{X} \times \mathcal{Y}} c_i^n d\gamma_i$  is bilinear, hence convex and concave in its variables, and continuous with respect to  $\lambda$ . Moreover, let  $(c_i^n)_n$  be non-decreasing sequences of bounded cost functions such that  $c_i = \sup_n c_i^n$ . By monotone convergence, we get  $f(\lambda, \gamma) = \sup_n \sum_{i=1}^N \lambda_i \int c_i^n d\gamma_i$ ,  $f(\lambda, \cdot)$ . So  $f$  the supremum of continuous functions, then  $f$  is lower semi-continuous with respect to  $\gamma$ , therefore Sion's minimax theorem (Sion, 1958) holds.

We are now able to prove Theorem 1.

**Proof.** Let  $\mathcal{X}$  and  $\mathcal{Y}$  be two Polish spaces. For all  $i \in \{1, \dots, N\}$ , we define  $c_i : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  a bounded below lower-semi cost function. The proof follows the exact same steps as those in the proof of (Villani, 2003, Theorem 1.3). First we suppose that  $\mathcal{X}$  and  $\mathcal{Y}$  are compact and that for all  $i$ ,  $c_i$  is continuous, then we show that it can be extended to  $\mathcal{X}$  and  $\mathcal{Y}$  non compact and finally to  $c_i$  only lower semi continuous.

First, let assume  $\mathcal{X}$  and  $\mathcal{Y}$  are compact and that for all  $i$ ,  $c_i$  is continuous. Let fix  $\lambda \in \Delta_N^+$ . We recall the topological dual of the space of bounded continuous functions  $\mathcal{C}^b(\mathcal{X} \times \mathcal{Y})$  endowed with  $\|\cdot\|_\infty$  norm, is the space of Radon measures  $\mathcal{M}(\mathcal{X} \times \mathcal{Y})$  endowed with total variation norm. We define, for  $u \in \mathcal{C}^b(\mathcal{X} \times \mathcal{Y})$ :

$$V_i^\lambda(u) = \begin{cases} 0 & \text{if } u \geq -\lambda_i c_i \\ +\infty & \text{else} \end{cases}$$

and:

$$E(u) = \begin{cases} \int f d\mu + \int g d\nu & \text{if } \exists (f, g) \in \mathcal{C}^b(\mathcal{X}) \times \mathcal{C}^b(\mathcal{Y}), u = f + g \\ +\infty & \text{else} \end{cases}$$

One can show that for all  $i$ ,  $V_i^\lambda$  is convex and lower semi-continuous (as the sublevel sets are closed) and  $E^\lambda$  is convex. More over for all  $i$ , these functions continuous in  $u_0 \equiv 1$  the hypothesis of Lemma 2 are satisfied.

Let now compute the Fenchel-Legendre transform of these function. Let  $\gamma \in \mathcal{M}(\mathcal{X} \times \mathcal{Y})$  :

$$\begin{aligned} V_i^{\lambda*}(-\gamma) &= \sup_{u \in \mathcal{C}^b(\mathcal{X} \times \mathcal{Y})} \left\{ - \int u d\gamma; \quad u \geq -\lambda_i c_i \right\} \\ &= \begin{cases} \int \lambda_i c_i d\gamma & \text{if } \gamma \in \mathcal{M}_+(\mathcal{X} \times \mathcal{Y}) \\ +\infty & \text{otherwise} \end{cases} \end{aligned}$$

On the other hand:

$$E^{\lambda*}(\gamma) = \begin{cases} 0 & \text{if } \forall (f, g) \in \mathcal{C}^b(\mathcal{X}) \times \mathcal{C}^b(\mathcal{Y}), \int f d\mu + \int g d\nu = \int (f + g) d\gamma \\ +\infty & \text{else} \end{cases}$$

This dual function is finite and equals 0 if and only if that the marginals of the dual variable  $\gamma$  are  $\mu$  and  $\nu$ .

Applying Lemma 2, we get:

$$\inf_{u \in \mathcal{C}^b(\mathcal{X} \times \mathcal{Y})} \sum_i V_i^\lambda(u) + E(u) = \sup_{\substack{\gamma_1, \dots, \gamma_N, \gamma \in \mathcal{M}(\mathcal{X} \times \mathcal{Y}) \\ \sum \gamma_i = \gamma}} \sum -V_i^{\lambda*}(\gamma_i) - E^{\lambda*}(-\gamma)$$

Hence, we have shown that, when  $\mathcal{X}$  and  $\mathcal{Y}$  are compact sets, and the costs  $(c_i)_i$  are continuous:

$$\sup_{(f, g) \in \mathcal{F}_c^\lambda} \int f d\mu + \int g d\nu = \inf_{\gamma \in \Gamma_{\mu, \nu}^N} \sum_i \lambda_i \int c_i d\gamma_i$$

Let now prove the result holds when the spaces  $\mathcal{X}$  and  $\mathcal{Y}$  are not compact. We still suppose that for all  $i$ ,  $c_i$  is uniformly continuous and bounded. We denote  $\|\mathbf{c}\|_\infty := \sup_i \sup_{(x, y) \in \mathcal{X} \times \mathcal{Y}} |c_i(x, y)|$ . Let define  $I^\lambda(\gamma) := \sum_i \lambda_i \int_{\mathcal{X} \times \mathcal{Y}} c_i d\gamma_i$

Let  $\gamma^* \in \Gamma_{\mu, \nu}^N$  such that  $I^\lambda(\gamma^*) = \min_{\gamma \in \Gamma_{\mu, \nu}^N} I^\lambda(\gamma)$ . The existence of the minimum comes from the lower-semi continuity of  $I^\lambda$  and the compacity of  $\Gamma_{\mu, \nu}^N$  for weak topology.

Let fix  $\delta \in (0, 1)$ .  $\mathcal{X}$  and  $\mathcal{Y}$  are Polish spaces then  $\exists \mathcal{X}_0 \subset \mathcal{X}, \mathcal{Y}_0 \subset \mathcal{Y}$  compacts such that  $\mu(\mathcal{X}_0^c) \leq \delta$  and  $\mu(\mathcal{Y}_0^c) \leq \delta$ . It follows that  $\forall i, \gamma_i^*((\mathcal{X}_0 \times \mathcal{Y}_0)^c) \leq 2\delta$ . Let define  $\gamma^{*0}$  such that for all  $i$ ,  $\gamma_i^{*0} = \frac{\mathbf{1}_{\mathcal{X}_0 \times \mathcal{Y}_0}}{\sum_i \gamma_i^*(\mathcal{X}_0 \times \mathcal{Y}_0)} \gamma_i^*$ . We define  $\mu_0 = \Pi_{1\#} \sum_i \gamma_i^{*0}$  and  $\nu_0 = \Pi_{2\#} \sum_i \gamma_i^{*0}$ . We then naturally define  $\Gamma_{0, \mu_0, \nu_0}^N := \{(\gamma_i)_{1 \leq i \leq N} \in \mathcal{M}_+(\mathcal{X}_0 \times \mathcal{Y}_0)^N \text{ s.t. } \Pi_{1\#} \sum_i \gamma_i = \mu_0 \text{ and } \Pi_{2\#} \sum_i \gamma_i = \nu_0\}$  and  $I_0^\lambda(\gamma_0) := \sum_i \lambda_i \int_{\mathcal{X}_0 \times \mathcal{Y}_0} c_i d\gamma_{0,i}$  for  $\gamma_0 \in \Gamma_{0, \mu_0, \nu_0}^N$ .

Let  $\tilde{\gamma}_0$  verifying  $I_0^\lambda(\tilde{\gamma}_0) = \min_{\gamma_0 \in \Gamma_{0, \mu_0, \nu_0}^N} I_0^\lambda(\gamma_0)$ . Let  $\tilde{\gamma} = (\sum_i \gamma_i^*(\mathcal{X}_0 \times \mathcal{Y}_0)) \tilde{\gamma}_0 + \mathbf{1}_{(\mathcal{X}_0 \times \mathcal{Y}_0)^c} \gamma^* \in \Gamma_{\mu, \nu}^N$ . Then we get

$$I^\lambda(\tilde{\gamma}) \leq \min_{\gamma_0 \in \Gamma_{0, \mu_0, \nu_0}^N} I_0^\lambda(\gamma_0) + 2 \sum |\lambda_i| \|\mathbf{c}\|_\infty \delta$$

We have already proved that:

$$\sup_{(f, g) \in \mathcal{F}_{0, \mathbf{c}}^\lambda} J_0^\lambda(f, g) = \inf_{\gamma_0 \in \Gamma_{0, \mu_0, \nu_0}^N} I_0^\lambda(\gamma_0)$$

with  $J_0^\lambda(f, g) = \int f d\mu_0 + \int g d\nu_0$  and  $\mathcal{F}_{0, \mathbf{c}}^\lambda$  is the set of  $(f, g) \in \mathcal{C}^b(\mathcal{X}_0) \times \mathcal{C}^b(\mathcal{Y}_0)$  satisfying, for every  $i$ ,  $f \oplus g \leq \min_i \lambda_i c_i$ . Let  $(\tilde{f}_0, \tilde{g}_0) \in \mathcal{F}_{0, \mathbf{c}}^\lambda$  such that :

$$J_0^\lambda(\tilde{f}_0, \tilde{g}_0) \geq \sup_{(f, g) \in \mathcal{F}_{0, \mathbf{c}}^\lambda} J_0^\lambda(f, g) - \delta$$

Since  $J_0^\lambda(0, 0) = 0$ , we get  $\sup J_0^\lambda \geq 0$  and then,  $J_0^\lambda(\tilde{f}_0, \tilde{g}_0) \geq \delta \geq -1$ . For every  $\gamma_0 \in \Gamma_{0, \mu_0, \nu_0}^N$ :

$$J_0^\lambda(\tilde{f}_0, \tilde{g}_0) = \int (\tilde{f}_0(x) + \tilde{g}_0(y)) d\gamma_0(x, y)$$

then we have the existence of  $(x_0, y_0) \in \mathcal{X}_0 \times \mathcal{Y}_0$  such that :  $\tilde{f}_0(x_0) + \tilde{g}_0(y_0) \geq -1$ . If we replace  $(\tilde{f}_0, \tilde{g}_0)$  by  $(\tilde{f}_0 - s, \tilde{g}_0 + s)$  for an accurate  $s$ , we get that:  $\tilde{f}_0(x_0) \geq \frac{1}{2}$  and  $\tilde{g}_0(y_0) \geq \frac{1}{2}$ , and then  $\forall (x, y) \in \mathcal{X}_0 \times \mathcal{Y}_0$ :

$$\begin{aligned} \tilde{f}_0(x) &\leq c'(x, y_0) - \tilde{g}_0(y_0) \leq c'(x, y_0) + \frac{1}{2} \\ \tilde{g}_0(y) &\leq c'(x_0, y) - \tilde{f}_0(x_0) \leq c'(x_0, y) + \frac{1}{2} \end{aligned}$$



where  $c' := \min_i \lambda_i c_i$ . Let define  $\bar{f}_0(x) = \inf_{y \in \mathcal{Y}_0} c'(x, y) - \bar{g}_0(y)$  for  $x \in \mathcal{X}$ . Then  $\bar{f}_0 \leq \bar{f}_0$  on  $\mathcal{X}_0$ . We then get  $J_0^\lambda(\bar{f}_0, \bar{g}_0) \geq J_0^\lambda(\tilde{f}_0, \tilde{g}_0)$  and  $\bar{f}_0 \leq c'(\cdot, y_0) + \frac{1}{2}$  on  $\mathcal{X}$ . Let define  $\bar{g}_0(y) = \inf_{x \in \mathcal{X}} c'(x, y) - \bar{f}_0(y)$ . By construction  $(f_0, g_0) \in \mathcal{F}_c^\lambda$  since the costs are uniformly continuous and bounded and  $J_0^\lambda(\bar{f}_0, \bar{g}_0) \geq J_0^\lambda(\tilde{f}_0, \tilde{g}_0) \geq J_0^\lambda(f_0, g_0)$ . We also have  $\bar{g}_0 \geq c'(x_0, \cdot) + \frac{1}{2}$  on  $\mathcal{Y}$ . Then we have in particular:  $\bar{g}_0 \geq -\|\mathbf{c}\|_\infty - \frac{1}{2}$  on  $\mathcal{X}$  and  $\bar{f}_0 \geq -\|\mathbf{c}\|_\infty - \frac{1}{2}$  on  $\mathcal{Y}$ . Finally:

$$\begin{aligned}
 J^\lambda(\bar{f}_0, \bar{g}_0) &:= \int_{\mathcal{X}_0} \bar{f}_0 d\mu_0 + \int_{\mathcal{Y}_0} \bar{g}_0 d\nu \\
 &= \sum_i \gamma_i^*(\mathcal{X}_0 \times \mathcal{Y}_0) \int_{\mathcal{X}_0 \times \mathcal{Y}_0} (\bar{f}_0(x) + \bar{g}_0(y)) d \left( \sum_i \gamma_i^{*0}(x, y) \right) \\
 &\quad + \int_{(\mathcal{X}_0 \times \mathcal{Y}_0)^c} \bar{f}_0(x) + \bar{g}_0(y) d \left( \sum_i \gamma_i^*(x, y) \right) \\
 &\geq (1 - 2\delta) \left( \int_{\mathcal{X}_0} \bar{f}_0 d\mu_0 + \int_{\mathcal{Y}_0} \bar{g}_0 d\nu_0 \right) - (2\|\mathbf{c}\|_\infty + 1) \sum_i \gamma_i^*((\mathcal{X}_0 \times \mathcal{Y}_0)^c) \\
 &\geq (1 - 2\delta) J_0^\lambda(\bar{f}_0, \bar{g}_0) - 2 \sum |\lambda_i| (2\|\mathbf{c}\|_\infty + 1) \delta \\
 &\geq (1 - 2\delta) J_0^\lambda(\tilde{f}_0, \tilde{g}_0) - 2 \sum |\lambda_i| (2\|\mathbf{c}\|_\infty + 1) \delta \\
 &\geq (1 - 2\delta) (\inf I_0^\lambda - \delta) - 2 \sum |\lambda_i| (2\|\mathbf{c}\|_\infty + 1) \delta \\
 &\geq (1 - 2\delta) (\inf I^\lambda - (2 \sum |\lambda_i| \|\mathbf{c}\|_\infty + 1) \delta) - 2 \sum |\lambda_i| (2\|\mathbf{c}\|_\infty + 1) \delta
 \end{aligned}$$

This being true for arbitrary small  $\delta$ , we get  $\sup J^\lambda \geq \inf I^\lambda$ . The other sens is always true then:

$$\sup_{(f,g) \in \mathcal{F}_c^\lambda} \int f d\mu + \int g d\nu = \inf_{\gamma \in \Gamma_{\mu,\nu}^N} \sum_i \lambda_i \int c_i d\gamma_i$$

for  $c_i$  uniformly continuous and  $\mathcal{X}$  and  $\mathcal{Y}$  non necessarily compact.

Let now prove that the result holds for lower semi-continuous costs. Let  $\mathbf{c} := (c_i)_i$  be a collection of lower semi-continuous costs. Let  $(c_i^n)_n$  be non-decreasing sequences of bounded below cost functions such that  $c_i = \sup_n c_i^n$ . Let fix  $\lambda \in \Delta_N^+$ . From last step, we have shown that for all  $n$ :

$$\inf_{\gamma \in \Gamma_{\mu,\nu}^N} I_n^\lambda(\gamma) = \sup_{(f,g) \in \mathcal{F}_{c^n}^\lambda} \int f d\mu + \int g d\nu \quad (18)$$

where  $I_n^\lambda(\gamma) = \sum_i \lambda_i \int c_i^n d\gamma_i$ . First it is clear that:

$$\sup_{(f,g) \in \mathcal{F}_c^\lambda} \int f d\mu + \int g d\nu \leq \sup_{(f,g) \in \mathcal{F}_{c^n}^\lambda} \int f d\mu + \int g d\nu \quad (19)$$

Let show that:

$$\inf_{\gamma \in \Gamma_{\mu,\nu}^N} I^\lambda(\gamma) = \sup_n \inf_{\gamma \in \Gamma_{\mu,\nu}^N} I_n^\lambda(\gamma) = \lim_n \inf_{\gamma \in \Gamma_{\mu,\nu}^N} I_n^\lambda(\gamma)$$

where  $I^\lambda(\gamma) = \sum_i \lambda_i \int c_i d\gamma_i$ .

Let  $(\gamma^{n,k})_k$  a minimizing sequence of  $\Gamma_{\mu,\nu}^N$  for the problem  $\inf_{\gamma \in \Gamma_{\mu,\nu}^N} \sum_i \lambda_i \int c_i^n d\gamma_i$ . By Lemma 1, up to an extraction, there exists  $\gamma^n \in \Gamma_{\mu,\nu}^N$  such that  $(\gamma^{n,k})_k$  converges weakly to  $\gamma^n$ . Then:

$$\inf_{\gamma \in \Gamma_{\mu,\nu}^N} I_n^\lambda(\gamma) = I_n^\lambda(\gamma^n)$$

Up to an extraction, there also exists  $\gamma^* \in \Gamma_{\mu,\nu}^N$  such that  $\gamma^n$  converges weakly to  $\gamma^*$ . For  $n \geq m$ ,  $I_n^\lambda(\gamma^n) \geq I_m^\lambda(\gamma^n) \geq I_m^\lambda(\gamma^m)$ , so by continuity of  $I_m^\lambda$ :

$$\lim_n I_n^\lambda(\gamma^n) \geq \limsup_n I_m^\lambda(\gamma^n) \geq I_m^\lambda(\gamma^*)$$

By monotone convergence,  $I_m^\lambda(\gamma^*) \rightarrow I^\lambda(\gamma^*)$  and  $\lim_n I_n^\lambda(\gamma^n) \geq I^\lambda(\gamma^*) \geq \inf_{\gamma \in \Gamma_{\mu,\nu}^N} I^\lambda(\gamma)$ .

Along with Eqs. 18 and 19, we get that:

$$\inf_{\gamma \in \Gamma_{\mu,\nu}^N} I^\lambda(\gamma) \leq \sup_{(f,g) \in \mathcal{F}_c^\lambda} \int f d\mu + \int g d\nu$$

The other sens being always true, we have then shown that, in the general case we still have:

$$\inf_{\gamma \in \Gamma_{\mu,\nu}^N} I^\lambda(\gamma) = \sup_{(f,g) \in \mathcal{F}_c^\lambda} \int f d\mu + \int g d\nu$$

To conclude, we apply Lemma 3, and we get:

$$\begin{aligned} \sup_{\lambda \in \Delta_N^+} \sup_{(f,g) \in \mathcal{F}_c^\lambda} \int f d\mu + \int g d\nu &= \sup_{\lambda \in \Delta_N^+} \inf_{\gamma \in \Gamma_{\mu,\nu}^N} I^\lambda(\gamma) \\ &= \text{EOT}_{\mathbf{c}}(\mu, \nu) \end{aligned}$$

## A.5 Proof of Proposition 3

**Proof 1.** Let recall that, from standard optimal transport results:

$$\text{EOT}_{\mathbf{c}}(\mu, \nu) = \sup_{u \in \Phi_{\mathbf{c}}} \int u d\mu d\nu$$

with  $\Phi_{\mathbf{c}} := \{u \in \mathcal{C}^b(\mathcal{X} \times \mathcal{Y}) \text{ s.t. } \exists \lambda \in \Delta_N^+, \exists \phi \in \mathcal{C}^b(\mathcal{X}), u = \phi^{cc} \oplus \phi^c \text{ with } c = \min_i \lambda_i c_i\}$  where  $\phi^c$  is the  $c$ -transform of  $\phi$ , i.e. for  $y \in \mathcal{Y}$ ,  $\phi^c(y) = \inf_{x \in \mathcal{X}} c(x, y) - \phi(x)$ .

Let denote  $\omega_1, \dots, \omega_N$  the continuity moduli of  $c_1, \dots, c_N$ . The existence of continuity moduli is ensured by the uniform continuity of  $c_1, \dots, c_N$  on the compact sets  $\mathcal{X} \times \mathcal{Y}$  (Heine's theorem). Then a modulus of continuity for  $\min_i \lambda_i c_i$  is  $\sum_i \lambda_i \omega_i$ . As  $\phi^c$  and  $\phi^{cc}$  share the same modulus of continuity than  $c = \min_i \lambda_i c_i$ , for  $u$  is  $\Phi_{\mathbf{c}}$ , a common modulus of continuity is  $2 \times \sum_i \omega_i$ . More over, it is clear that for all  $x, y$ ,  $\{u(x, y) \text{ s.t. } u \in \Phi_{\mathbf{c}}\}$  is compact. Then, applying Ascoli's theorem, we get, that  $\Phi_{\mathbf{c}}$  is compact for  $\|\cdot\|_\infty$  norm. By continuity of  $u \rightarrow \int u d\mu d\nu$ , the supremum is attained, and we get the existence of the optimum  $u^*$ . The existence of optima  $(\lambda^*, f^*, g^*)$  immediately follows.

Let first assume that  $(\gamma_k)_{k=1}^N$  is a solution of Eq. (1) and  $(\lambda, f, g)$  is a solution of Eq. (5). Then it is clear that for all  $i, j$ ,  $f \oplus g \leq \lambda_i c_i$ ,  $(\gamma_k)_{k=1}^N \in \Gamma_{\mu,\nu}^N$  and  $\int c_j d\gamma_j = \int c_i d\gamma_i$  (by Proposition 1). Let  $k \in \{1, \dots, N\}$ . Moreover, by Theorem 1:

$$\begin{aligned} 0 &= \int f d\mu + \int g d\nu - \int c_i d\gamma_i \\ &= \sum \int (f(x) + g(y)) d\gamma_i(x, y) - \sum_i \lambda_i \int c_i(x, y) d\gamma_i(x, y) \\ &= \sum \int (f(x) + g(y) - \lambda_i c_i(x, y)) d\gamma_i(x, y) \end{aligned}$$

Since  $f \oplus g \leq \lambda_i c_i$  and  $\gamma_i$  are positive measures then  $f \oplus g = \lambda_i c_i$ ,  $\gamma_i$ -almost everywhere.

Reciprocally, let assume that there exist  $(\gamma_k)_{k=1}^N \in \Gamma_{\mu, \nu}^N$  and  $(\lambda, f, g) \in \Delta_n^+ \times \mathcal{C}^b(\mathcal{X}) \times \mathcal{C}^b(\mathcal{Y})$  such that  $\forall i \in \{1, \dots, N\}$ ,  $f \oplus g \leq \lambda_i c_i$ ,  $\forall i, j \in \{1, \dots, N\}$   $\int c_i d\gamma_i = \int c_j d\gamma_j$  and  $f \oplus g = \lambda_i c_i$   $\gamma_i$ -a.e.. Then, for any  $k$ :

$$\begin{aligned} \int c_k d\gamma_k &= \sum_i \lambda_i \int c_i d\gamma_i \\ &= \sum_i \int (f(x) + g(y)) d\gamma_i(x, y) \\ &= \int f(x) d\mu(x) + \int g(y) d\nu(y) \\ &\leq \text{EOT}_{\mathbf{c}}(\mu, \nu) \text{ by Theorem 1} \end{aligned}$$

then  $\gamma_k$  is solution of the primal problem. We also have for any  $k$ :

$$\begin{aligned} \int f d\mu + \int g d\nu &= \sum_i \int (f(x) + g(y)) d\gamma_i(x, y) \\ &= \sum_i \int \lambda_i c_i d\gamma_i \\ &= \int c_k d\gamma_k \\ &\geq \text{EOT}_{\mathbf{c}}(\mu, \nu) \end{aligned}$$

then, thanks to Theorem 1,  $(\lambda, f, g)$  is solution of the dual problem.

Let now proof the result stated in Remark 2. Let assume the costs are strictly positive or strictly negative. If there exist  $i$  such that  $\lambda_i = 0$ , thanks to the condition  $f \oplus g \leq \lambda_i c_i$ , we get  $f \oplus g \leq 0$  and then  $f \oplus g = 0$  which contradicts the conditions  $f \oplus g = \lambda_k c_k$  for all  $k$ .

## A.6 Proof of Proposition 4

Before proving the result let us first introduce the following lemma.

**Lemma 4.** Let  $\mathcal{X}$  and  $\mathcal{Y}$  be Polish spaces. Let  $\mathbf{c} := (c_i)_{1 \leq i \leq N}$  a family of bounded below continuous costs. For  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  and  $\lambda \in \Delta_N^+$ , we define

$$c_\lambda(x, y) := \min_{i=1, \dots, N} (\lambda_i c_i(x, y))$$

then for any  $(\mu, \nu) \in \mathcal{M}_+^1(\mathcal{X}) \times \mathcal{M}_+^1(\mathcal{Y})$

$$\text{EOT}_{\mathbf{c}}(\mu, \nu) = \sup_{\lambda \in \Delta_N^+} W_{c_\lambda}(\mu, \nu) \quad (20)$$

**Proof.** Let  $(\mu, \nu) \in \mathcal{M}_+^1(\mathcal{X}) \times \mathcal{M}_+^1(\mathcal{Y})$  and  $\mathbf{c} := (c_i)_{1 \leq i \leq N}$  cost functions on  $\mathcal{X} \times \mathcal{Y}$ . Let  $\lambda \in \Delta_N^+$ , then by Proposition 1:

$$\text{EOT}_{\mathbf{c}}(\mu, \nu) = \sup_{\lambda \in \Delta_N^+} \sup_{(f, g) \in \mathcal{F}_{\mathbf{c}}^\lambda} \int_{\mathcal{X}} f(x) d\mu(x) + \int_{\mathcal{Y}} g(y) d\nu(y)$$

Therefore by denoting  $c_\lambda := \min_i (\lambda_i c_i)$  which is a continuous. The dual form of the classical Optimal Transport problem gives that:

$$\sup_{(f, g) \in \mathcal{F}_{\mathbf{c}}^\lambda} \int_{\mathcal{X}} f(x) d\mu(x) + \int_{\mathcal{Y}} g(y) d\nu(y) = W_{c_\lambda}(\mu, \nu)$$

and the result follows.

Let us now prove the result of Proposition 4.

**Proof.** Let  $\mu$  and  $\nu$  be two probability measures. Let  $\alpha \in (0, 1]$ . Note that if  $d$  is a metric then  $d^\alpha$  too. Therefore in the following we consider  $d$  a general metric on  $\mathcal{X} \times \mathcal{X}$ . Let  $c_1 : (x, y) \rightarrow 2 \times \mathbf{1}_{x \neq y}$  and  $c_2 = d^\alpha$ . For all  $\lambda \in [0, 1]$ :

$$c_\lambda(x, y) := \min(\lambda c_1(x, y), (1 - \lambda) c_2(x, y)) = \min(2\lambda, (1 - \lambda) d(x, y))$$

defines a distance on  $\mathcal{X} \times \mathcal{X}$ . Then according to (Villani, 2003, Theorem 1.14):

$$W_{c_\lambda}(\mu, \nu) = \sup_{f \text{ s.t. } f \text{ } 1-c_\lambda \text{ Lipschitz}} \int f d\mu - \int f d\nu$$

Then thanks to Lemma 4 we have

$$\text{EOT}_{(c_1, c_2)}(\mu, \nu) = \sup_{\lambda \in [0, 1], f \text{ s.t. } f \text{ } 1-c_\lambda \text{ Lipschitz}} \int f d\mu - \int f d\nu$$

Let now prove that in this case:  $\text{EOT}_{(c_1, c_2)}(\mu, \nu) = \beta_d(\mu, \nu)$ . Let  $\lambda \in [0, 1]$  and  $f$  a  $c_\lambda$  Lipschitz function.  $f$  is lower bounded: let  $m = \inf f$  and  $(u_n)_n$  a sequence satisfying  $f(u_n) \rightarrow m$ . Then for all  $x, y$ ,  $f(x) - f(y) \leq 2\lambda$  and  $f(x) - f(y) \leq (1 - \lambda)d(x, y)$ . Let define  $g = f - m - \lambda$ . For  $x$  fixed and for all  $n$ ,  $f(x) - f(u_n) \leq 2\lambda$ , so taking the limit in  $n$  we get  $f(x) - m \leq 2\lambda$ . So we get that for all  $x, y$ ,  $g(x) \in [-\lambda, +\lambda]$  and  $g(x) - g(y) \in [-(1 - \lambda)d(x, y), (1 - \lambda)d(x, y)]$ . Then  $\|g\|_\infty \leq \lambda$  and  $\|g\|_d \leq 1 - \lambda$ . By construction, we also have  $\int f d\mu - \int f d\nu = \int g d\mu - \int g d\nu$ . Then  $\|g\|_\infty + \|g\|_d \leq 1$ . So we get that  $\text{EOT}_{(c_1, c_2)}(\mu, \nu) \leq \beta_d(\mu, \nu)$ . Reciprocally, let  $g$  be a function satisfying  $\|g\|_\infty + \|g\|_d \leq 1$ . Let define  $f = g + \|g\|_\infty$  and  $\lambda = \|g\|_\infty$ . Then, for all  $x, y$ ,  $f(x) \in [0, 2\lambda]$  and so  $f(x) - f(y) \leq 2\lambda$ . It is immediate that  $f(x) - f(y) \in [-(1 - \lambda)d(x, y), (1 - \lambda)d(x, y)]$ . Then we get  $f(x) - f(y) \leq \min(\lambda, (1 - \lambda)d(x, y))$ . And by construction, we still have  $\int f d\mu - \int f d\nu = \int g d\mu - \int g d\nu$ . So  $\text{EOT}_{(c_1, c_2)}(\mu, \nu) \geq \beta_d(\mu, \nu)$ .

Finally we get  $\text{EOT}_{(c_1, c_2)}(\mu, \nu) = \beta_d(\mu, \nu)$  when  $c_1 : (x, y) \rightarrow 2 \times \mathbf{1}_{x \neq y}$  and  $c_2 = d$  a distance on  $\mathcal{X} \times \mathcal{X}$ .

## A.7 Proof of Proposition 5

**Lemma 5.** Let  $x_1, \dots, x_N \geq 0$ , then:

$$\sup_{\lambda \in \Delta_N^+} \min_i \lambda_i x_i = \frac{1}{\sum_i \frac{1}{x_i}}$$

**Proof.** First if there exists  $i$  such that  $x_i = 0$ , we immediately have  $\sup_{\lambda \in \Delta_N^+} \min_i \lambda_i x_i = 0$ .

$g : \lambda \mapsto \min_i \lambda_i x_i$  is a continuous function on the compact set  $\lambda \in \Delta_N^+$ . Let denote  $\lambda^*$  the maximum of  $g$ .

Let show that for all  $i, j$ ,  $\lambda_i^* x_i = \lambda_j^* x_j$ . Let denote  $i_0, \dots, i_k$  the indices such that  $\lambda_{i_l}^* x_{i_l} = \min_i \lambda_i^* x_i$ . Let assume there exists  $j_0$  such that:  $\lambda_{j_0}^* x_{j_0} > \min_i \lambda_i^* x_i$ , and that all other indices  $i$  have a larger  $\lambda_i^* x_i \geq \lambda_{j_0}^* x_{j_0}$ . Then for  $\epsilon > 0$  sufficiently small, let  $\tilde{\lambda}$  defined as:  $\tilde{\lambda}_{j_0} = \lambda_{j_0}^* - \epsilon$ ,  $\tilde{\lambda}_{i_l} = \lambda_{i_l}^* + \epsilon/k$  for all  $l \in \{1, \dots, k\}$  and  $\tilde{\lambda}_i = \lambda_i^*$  for all other indices. Then  $\tilde{\lambda} \in \Delta_N^+$  and  $g(\tilde{\lambda}) < g(\lambda^*)$ , which contradicts that  $\lambda^*$  is the maximum.

Then at the optimum for all  $i, j$ ,  $\lambda_i^* x_i = \lambda_j^* x_j$ . So  $\lambda_i^* x_i = C$  for a certain constant  $C$ . Moreover  $\sum_i \lambda_i^* = 1$ . Then  $1/C = \sum_i 1/x_i$ . Finally, for all  $i$ ,

$$\lambda_i^* = \frac{1/x_i}{\sum_i 1/x_i}$$

and then:

$$\sup_{\lambda \in \Delta_N^+} \min_i \lambda_i x_i = \frac{1}{\sum_i \frac{1}{x_i}}.$$

**Proof.** Let  $\mu$  and  $\nu$  be two probability measures respectively on  $\mathcal{X}$  and  $\mathcal{Y}$ . Let  $\mathbf{c} := (c_i)_i$  be a family of cost functions. Let define for  $\lambda \in \Delta_N^+$ ,  $c_\lambda(x, y) := \min_i (\lambda_i c_i(x, y))$ . We have, by linearity  $W_{c_\lambda}(\mu, \nu) \leq \min_i (\lambda_i W_{c_i}(\mu, \nu))$ . So we deduce by Lemma 4:

$$\begin{aligned}
 \text{EOT}_{\mathbf{c}}(\mu, \nu) &= \sup_{\lambda \in \Delta_N^+} W_{c_\lambda}(\mu, \nu) \\
 &\leq \sup_{\lambda \in \Delta_N^+} \min_i \lambda_i W_{c_i}(\mu, \nu) \\
 &= \frac{1}{\sum_i \frac{1}{W_{c_i}(\mu, \nu)}} \text{ by Lemma 5}
 \end{aligned}$$

which concludes the proof.

### A.8 Proof of Theorem 2

**Proof.** To show the strong duality of the regularized problem, we use the same sketch of proof as for the strong duality of the original problem. Let first assume that, for all  $i$ ,  $c_i$  is continuous on the compact set  $\mathcal{X} \times \mathcal{Y}$ . Let fix  $\lambda \in \Delta_N^+$ . We define, for all  $u \in \mathcal{C}^b(\mathcal{X} \times \mathcal{Y})$ :

$$V_i^\lambda(u) = \varepsilon_i \left( \int_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \exp \frac{-u(x,y) - \lambda_i c_i(x,y)}{\varepsilon_i} d\mu(x) d\nu(y) - 1 \right)$$

and:

$$E(u) = \begin{cases} \int f d\mu + \int g d\nu & \text{if } \exists (f, g) \in \mathcal{C}^b(\mathcal{X}) \times \mathcal{C}^b(\mathcal{Y}), u = f + g \\ +\infty & \text{else} \end{cases}$$

Let compute the Fenchel-Legendre transform of these functions. Let  $\gamma \in \mathcal{M}(\mathcal{X} \times \mathcal{Y})$ :

$$V_i^{\lambda*}(-\gamma) = \sup_{u \in \mathcal{C}^b(\mathcal{X} \times \mathcal{Y})} - \int u d\gamma - \varepsilon_i \left( \int_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \exp \frac{-u(x,y) - \lambda_i c_i(x,y)}{\varepsilon_i} d\mu(x) d\nu(y) - 1 \right)$$

However, by density of  $\mathcal{C}^b(\mathcal{X} \times \mathcal{Y})$  in  $L^1_{d\mu \otimes \nu}(\mathcal{X} \times \mathcal{Y})$ , the set of integrable functions for  $\mu \otimes \nu$  measure, we deduce that

$$V_i^{\lambda*}(-\gamma) = \sup_{u \in L^1_{d\mu \otimes \nu}(\mathcal{X} \times \mathcal{Y})} - \int u d\gamma - \varepsilon_i \left( \int_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \exp \frac{-u(x,y) - \lambda_i c_i(x,y)}{\varepsilon_i} d\mu(x) d\nu(y) - 1 \right)$$

This supremum equals  $+\infty$  if  $\gamma$  is not positive and not absolutely continuous with regard to  $\mu \otimes \nu$ . Let us now denote  $F_{\gamma, \lambda}(u) := - \int u d\gamma - \varepsilon_i \left( \int_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \exp \frac{-u(x,y) - \lambda_i c_i(x,y)}{\varepsilon_i} d\mu(x) d\nu(y) - 1 \right)$ .  $F_{\gamma, \lambda}$  is Fréchet differentiable and its maximum is attained for  $u^* = \varepsilon_i \log \left( \frac{d\gamma}{d\mu \otimes \nu} \right) + \lambda_i c_i$ . Therefore we obtain that

$$\begin{aligned}
 V_i^{\lambda*}(-\gamma) &= \varepsilon_i \left( \int \log \left( \frac{d\gamma}{d\mu \otimes \nu} \right) d\gamma + 1 - \gamma(\mathcal{X} \times \mathcal{Y}) \right) + \lambda_i \int c_i d\gamma \\
 &= \lambda_i \int c_i d\gamma + \varepsilon_i \text{KL}(\gamma_i || \mu \times \nu)
 \end{aligned}$$

Thanks to the compactness of  $\mathcal{X} \times \mathcal{Y}$ , all the  $V_i^\lambda$  for  $i \in \{1, \dots, N\}$  are continuous on  $\mathcal{C}^b(\mathcal{X} \times \mathcal{Y})$ . Therefore by applying Lemma 2, we obtain that:

$$\inf_{u \in \mathcal{C}^b(\mathcal{X} \times \mathcal{Y})} \sum_i V_i^\lambda(u) + E(u) = \sup_{\gamma_1, \dots, \gamma_N, \gamma \in \mathcal{M}(\mathcal{X} \times \mathcal{Y})} - \sum_i V_i^{\lambda*}(\gamma_i) - E^*(-\gamma)$$

$$\begin{aligned}
 & \sup_{f \in \mathcal{C}^b(\mathcal{X}), g \in \mathcal{C}^b(\mathcal{Y})} \int f d\mu + \int g d\nu \\
 & - \sum_{i=1}^N \varepsilon_i \left( \int_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \exp \frac{f(x) + g(y) - \lambda_i c_i(x,y)}{\varepsilon_i} d\mu(x) d\nu(y) - 1 \right) \\
 & = \inf_{\gamma \in \Gamma_{\mu,\nu}^N} \sum_{i=1}^N \lambda_i \int c_i d\gamma_i + \varepsilon_i \text{KL}(\gamma_i || \mu \otimes \nu)
 \end{aligned}$$

Therefore by considering the supremum over the  $\lambda \in \Delta_N$ , we obtain that

$$\begin{aligned}
 & \sup_{\lambda \in \Delta_N^+} \sup_{f \in \mathcal{C}^b(\mathcal{X}), g \in \mathcal{C}^b(\mathcal{Y})} \int f d\mu + \int g d\nu \\
 & - \sum_{i=1}^N \varepsilon_i \left( \int_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \exp \frac{f(x) + g(y) - \lambda_i c_i(x,y)}{\varepsilon_i} d\mu(x) d\nu(y) - 1 \right) \\
 & = \sup_{\lambda \in \Delta_N^+} \inf_{\gamma \in \Gamma_{\mu,\nu}^N} \sum_{i=1}^N \lambda_i \int c_i d\gamma_i + \varepsilon_i \text{KL}(\gamma_i || \mu \otimes \nu)
 \end{aligned}$$

Let  $f : (\lambda, \gamma) \in \Delta_N^+ \times \Gamma_{\mu,\nu}^N \mapsto \sum_{i=1}^N \lambda_i \int c_i d\gamma_i + \varepsilon_i \text{KL}(\gamma_i || \mu \otimes \nu)$ .  $f$  is clearly concave and continuous in  $\lambda$ . Moreover  $\gamma \mapsto \text{KL}(\gamma_i || \mu \otimes \nu)$  is convex and lower semi-continuous for weak topology (Dupuis and Ellis, 2011, Lemma 1.4.3). Hence  $f$  is convex and lower-semi continuous in  $\gamma$ .  $\Delta_N^+$  is convex, and  $\Gamma_{\mu,\nu}^N$  is compact for weak topology (see Lemma 1). So by Sion's theorem, we get the expected result:

$$\begin{aligned}
 & \min_{\gamma \in \Gamma_{\mu,\nu}^N} \sup_{\lambda \in \Delta_N^+} \sum_i \lambda_i \int c_i d\gamma_i + \sum_i \varepsilon_i \text{KL}(\gamma_i || \mu \otimes \nu) \\
 & = \sup_{\lambda \in \Delta_N^+} \sup_{(f,g) \in \mathcal{C}_b(\mathcal{X}) \times \mathcal{C}_b(\mathcal{Y})} \int_{\mathcal{X}} f(x) d\mu(x) + \int_{\mathcal{Y}} g(y) d\nu(y) \\
 & - \sum_{i=1}^N \varepsilon_i \left( \int_{\mathcal{X} \times \mathcal{Y}} e^{\frac{f(x) + g(y) - \lambda_i c_i(x,y)}{\varepsilon_i}} d\mu(x) d\nu(y) - 1 \right)
 \end{aligned}$$

Moreover by fixing  $\gamma \in \Gamma_{\mu,\nu}^N$ , we have

$$\begin{aligned}
 & \sup_{\lambda \in \Delta_N^+} \sum_i \lambda_i \int c_i d\gamma_i + \sum_i \varepsilon_i \text{KL}(\gamma_i || \mu \otimes \nu) \\
 & = \max_i \int c_i d\gamma_i + \sum_j \varepsilon_j \text{KL}(\gamma_j || \mu \otimes \nu)
 \end{aligned}$$

which concludes the proof in case of continuous costs. A similar proof as the one of the Theorem 2 allows to extend the results for lower semi-continuous cost functions.

## B Discrete cases

### B.1 Exact discrete case

Let  $a \in \Delta_N^+$  and  $b \in \Delta_m^+$  and  $\mathbf{C} := (C_i)_{1 \leq i \leq N} \in (\mathbb{R}^{n \times m})^N$  be  $N$  cost matrices. Let also  $\mathbf{X} := \{x_1, \dots, x_n\}$  and  $\mathbf{Y} := \{y_1, \dots, y_m\}$  two subset of  $\mathcal{X}$  and  $\mathcal{Y}$  respectively. Moreover we define the two following discrete measure  $\mu = \sum_{i=1}^n a_i \delta_{x_i}$  and  $\nu = \sum_{i=1}^m b_i \delta_{y_i}$  and for all  $i$ ,  $C_i = (c_i(x_k, y_l))_{1 \leq k \leq n, 1 \leq l \leq m}$  where  $(c_i)_{i=1}^N$  a family of cost functions. The discretized multiple cost optimal transport primal problem can be written as follows:

$$\text{EOT}_c(\mu, \nu) = \widehat{\text{EOT}}_{\mathbf{C}}(a, b) := \inf_{P \in \Gamma_{a,b}^N} \max_i \langle P_i, C_i \rangle$$

where  $\Gamma_{a,b}^N := \left\{ (P_i)_{1 \leq i \leq N} \in (\mathbb{R}_+^{n \times m})^N \text{ s.t. } (\sum_i P_i) \mathbf{1}_m = a \text{ and } (\sum_i P_i^T) \mathbf{1}_n = b \right\}$ . As in the continuous case, strong duality holds and we can rewrite the dual in the discrete case also.

**Proposition 6** (Duality for the discrete problem). *Let  $a \in \Delta_N^+$  and  $b \in \Delta_m^+$  and  $\mathbf{C} := (C_i)_{1 \leq i \leq N} \in (\mathbb{R}^{n \times m})^N$  be  $N$  cost matrices. Strong duality holds for the discrete problem and*

$$\widehat{\text{EOT}}_{\mathbf{C}}(a, b) = \sup_{\lambda \in \Delta_N^+} \sup_{(f,g) \in \mathcal{F}_{\mathbf{C}}^\lambda} \langle f, a \rangle + \langle g, b \rangle.$$

where  $\mathcal{F}_{\mathbf{C}}^\lambda := \{(f, g) \in \mathbb{R}_+^n \times \mathbb{R}_+^m \text{ s.t. } \forall i \in \{1, \dots, N\}, f \mathbf{1}_m^T + \mathbf{1}_n g^T \leq \lambda_i C_i\}$ .

### B.2 Entropic regularized discrete case

We now extend the regularization in the discrete case. Let  $a \in \Delta_n^+$  and  $b \in \Delta_m^+$  and  $\mathbf{C} := (C_i)_{1 \leq i \leq N} \in (\mathbb{R}^{n \times m})^N$  be  $N$  cost matrices and  $\varepsilon = (\varepsilon_i)_{1 \leq i \leq N}$  be nonnegative real numbers. The discretized regularized primal problem is:

$$\widehat{\text{EOT}}_{\mathbf{C}}^\varepsilon(a, b) = \inf_{P \in \Gamma_{a,b}^N} \max_i \langle P_i, C_i \rangle - \sum_{i=1}^N \varepsilon_i H(P_i)$$

where  $H(P) = \sum_{i,j} P_{i,j} (\log P_{i,j} - 1)$  for  $P = (P_{i,j})_{i,j} \in \mathbb{R}_+^{n \times m}$  is the discrete entropy. In the discrete case, strong duality holds thanks to Lagrangian duality and Slater sufficient conditions:

**Proposition 7** (Duality for the discrete regularized problem). *Let  $a \in \Delta_n^+$  and  $b \in \Delta_m^+$  and  $\mathbf{C} := (C_i)_{1 \leq i \leq N} \in (\mathbb{R}^{n \times m})^N$  be  $N$  cost matrices and  $\varepsilon := (\varepsilon_i)_{1 \leq i \leq N}$  be non negative reals. Strong duality holds and by denoting  $K_i^{\lambda_i} = \exp(-\lambda_i C_i / \varepsilon_i)$ , we have*

$$\widehat{\text{EOT}}_{\mathbf{C}}^\varepsilon(a, b) = \sup_{\lambda \in \Delta_N^+} \sup_{f \in \mathbb{R}^n, g \in \mathbb{R}^m} \langle f, a \rangle + \langle g, b \rangle - \sum_{i=1}^N \varepsilon_i \langle e^{f/\varepsilon_i}, K_i^{\lambda_i} e^{g/\varepsilon_i} \rangle.$$

The objective function for the dual problem is strictly concave in  $(\lambda, f, g)$  but is neither smooth or strongly convex.

**Proof.** *The proofs in the discrete case are simpler and only involves Lagrangian duality (Boyd et al., 2004, Chapter 5). Let do the proof in the regularized case, the one for the standard problem follows exactly the same path.*



Let  $a \in \Delta_N^+$  and  $b \in \Delta_m^+$  and  $\mathbf{C} := (C_i)_{1 \leq i \leq N} \in (\mathbb{R}^{n \times m})^N$  be  $N$  cost matrices.

$$\begin{aligned}
 \widehat{\text{EOT}}_{\mathbf{C}}^{\varepsilon}(a, b) &= \inf_{P \in \Gamma_{a,b}^N} \max_{1 \leq i \leq N} \langle P_i, C_i \rangle - \sum_{i=1}^N \varepsilon_i \text{H}(P_i) \\
 &= \inf_{\substack{(t, P) \in \mathbb{R} \times (\mathbb{R}_+^{n \times m})^N \\ (\sum_i P_i) \mathbf{1}_m = a \\ (\sum_i P_i^T) \mathbf{1}_n = b \\ \forall j, \langle P_j, C_j \rangle \leq t}} t - \sum_{i=1}^N \varepsilon_i \text{H}(P_i) \\
 &= \inf_{(t, P) \in \mathbb{R} \times (\mathbb{R}_+^{n \times m})^N} \sup_{f \in \mathbb{R}^n, g \in \mathbb{R}^m, \lambda \in \mathbb{R}_+^N} t + \sum_{j=1}^N \lambda_j (\langle P_j, C_j \rangle - t) - \sum_{i=1}^N \varepsilon_i \text{H}(P_i) \\
 &\quad + f^T \left( a - \sum_i P_i \mathbf{1}_m \right) + g^T \left( b - \sum_i P_i^T \mathbf{1}_n \right)
 \end{aligned}$$

The constraints are qualified for this convex problem, hence by Slater's sufficient condition (Boyd et al., 2004, Section 5.2.3), strong duality holds and:

$$\begin{aligned}
 \widehat{\text{EOT}}_{\mathbf{C}}^{\varepsilon}(a, b) &= \sup_{f \in \mathbb{R}^n, g \in \mathbb{R}^m, \lambda \in \mathbb{R}_+^N} \inf_{(t, P) \in \mathbb{R} \times (\mathbb{R}_+^{n \times m})^N} t + \sum_{j=1}^N \lambda_j (\langle P_j, C_j \rangle - t) - \sum_{j=1}^N \varepsilon_j \text{H}(P_j) \\
 &\quad + f^T \left( a - \sum_{j=1}^N P_j \mathbf{1}_m \right) + g^T \left( b - \sum_{j=1}^N P_j^T \mathbf{1}_n \right) \\
 &= \sup_{\substack{f \in \mathbb{R}^n \\ g \in \mathbb{R}^m \\ \lambda \in \Delta_N^+}} \langle f, a \rangle + \langle g, b \rangle + \sum_{j=1}^N \inf_{P_j \in \mathbb{R}_+^{n \times m}} (\langle P_j, \lambda_j C_j - f \mathbf{1}_n^T - \mathbf{1}_m g^T \rangle - \varepsilon_j \text{H}(P_j))
 \end{aligned}$$

But for every  $i = 1, \dots, N$  the solution of

$$\inf_{P_j \in \mathbb{R}_+^{n \times m}} (\langle P_j, \lambda_j C_j - f \mathbf{1}_n^T - \mathbf{1}_m g^T \rangle - \varepsilon_j \text{H}(P_j))$$

is

$$P_j = \exp \left( \frac{f \mathbf{1}_n^T + \mathbf{1}_m g^T - \lambda_j C_j}{\varepsilon_j} \right)$$

Finally we obtain that

$$\widehat{\text{EOT}}_{\mathbf{C}}^{\varepsilon}(a, b) = \sup_{f \in \mathbb{R}^n, g \in \mathbb{R}^m, \lambda \in \Delta_N^+} \langle f, a \rangle + \langle g, b \rangle - \sum_{k=1}^N \varepsilon_k \sum_{i,j} \exp \left( \frac{f_i + g_j - \lambda_k C_k^{i,j}}{\varepsilon_k} \right)$$

## C Other results

### C.1 Utilitarian and Optimal Transport

**Proposition 8.** *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be Polish spaces. Let  $\mathbf{c} := (c_i)_{1 \leq i \leq N}$  be a family of bounded below continuous cost functions on  $\mathcal{X} \times \mathcal{Y}$ , and  $\mu \in \mathcal{M}_+^1(\mathcal{X})$  and  $\nu \in \mathcal{M}_+^1(\mathcal{Y})$ . Then we have:*

$$\inf_{(\gamma_i)_{i=1}^N \in \Gamma_{\mu, \nu}^N} \sum_i \int c_i d\gamma_i = W_{\min_i(c_i)}(\mu, \nu) \quad (21)$$

**Proof.** *The proof is a by-product of the proof of Theorem 1. The continuity of the costs is necessary since  $\min_i(c_i)$  is not necessarily lower semi-continuous when the costs are supposed lower semi-continuous.*

**Remark 3.** *We thank an anonymous reviewer for noticing that the utilitarian problem can be written also as an Optimal Transport on the space  $\mathcal{Z} = (\mathcal{X} \times \{1, \dots, N\}) \times (\mathcal{Y} \times \{1, \dots, N\})$ :*

$$\min_{\gamma \in \tilde{\Gamma}_{\mu, \nu}} \int_{x, i, y, j} c((x, i), (y, j)) d\gamma(x, i, y, j)$$

where the constraint space is  $\tilde{\Gamma}_{\mu, \nu} := \{\gamma \in \mathcal{M}_+^1(\mathcal{Z}) \text{ s.t. } \Pi_{\mathcal{X}}\gamma = \mu, \Pi_{\mathcal{Y}}\gamma = \nu\}$ .

### C.2 MOT generalizes OT

**Proposition 9.** *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be Polish spaces. Let  $N \geq 0$ ,  $\mathbf{c} = (c_i)_{1 \leq i \leq N}$  be a family of nonnegative lower semi-continuous costs and let us denote for all  $k \in \{1, \dots, N\}$ ,  $\mathbf{c}_k = (c_i)_{1 \leq i \leq k}$ . Then for all  $k \in \{1, \dots, N\}$ , there exists a family of costs  $\mathbf{d}_k \in LSC(\mathcal{X} \times \mathcal{Y})^N$  such that*

$$\text{EOT}_{\mathbf{d}_k}(\mu, \nu) = \text{EOT}_{\mathbf{c}_k}(\mu, \nu) \quad (22)$$

**Proof.** *For all  $k \in \{1, \dots, N\}$ , we define  $\mathbf{d}_k := (c_1, \dots, (N - k + 1) \times c_k, \dots, (N - k + 1) \times c_k)$ . Therefore, thanks to Lemma 4 we have*

$$\text{EOT}_{\mathbf{d}_k}(\mu, \nu) = \sup_{\lambda \in \Delta_N^+} W_{c_\lambda}(\mu, \nu) \quad (23)$$

$$= \sup_{(\lambda, \gamma) \in \Delta_n^k} \inf_{\gamma \in \Gamma_{\mu, \nu}} \int_{\mathcal{X} \times \mathcal{Y}} \min(\lambda_1 c_1, \dots, \lambda_{k-1} c_{k-1}, \lambda_k c_k) d\gamma \quad (24)$$

where  $\Delta_n^k := \{(\lambda, \gamma) \in \Delta_N^+ \times \mathbb{R}_+ : \gamma = (N - k + 1) \times \min(\lambda_k, \dots, \lambda_N)\}$ . First remarks that

$$\gamma = 1 - \sum_{i=1}^{k-1} \lambda_i \iff (N - k + 1) \times \min(\lambda_k, \dots, \lambda_N) = \sum_{i=k}^N \lambda_i \quad (25)$$

$$\iff \lambda_k = \dots = \lambda_N \quad (26)$$

But in that case  $(\lambda_1, \dots, \lambda_{k-1}, \gamma) \in \Delta_k$  and therefore we obtain that

$$\text{EOT}_{\mathbf{d}_k}(\mu, \nu) \geq \sup_{\lambda \in \Delta_k} \inf_{\gamma \in \Gamma_{\mu, \nu}} \int_{\mathcal{X} \times \mathcal{Y}} \min(\lambda_1 c_1, \dots, \lambda_{k-1} c_{k-1}, \gamma c_k) d\gamma = \text{EOT}_{\mathbf{c}_k}(\mu, \nu)$$

Finally by definition we have  $\gamma \leq \sum_{i=k}^N \lambda_i = 1 - \sum_{i=1}^{k-1} \lambda_i$  and therefore

$$\int_{\mathcal{X} \times \mathcal{Y}} \min(\lambda_1 c_1, \dots, \lambda_{k-1} c_{k-1}, \gamma c_k) d\gamma \leq \int_{\mathcal{X} \times \mathcal{Y}} \min \left( \lambda_1 c_1, \dots, \lambda_{k-1} c_{k-1}, \left(1 - \sum_{i=1}^{k-1} \lambda_i\right) c_k \right)$$

Then we obtain that

$$\text{EOT}_{\mathbf{d}_k}(\mu, \nu) \leq \text{EOT}_{\mathbf{c}_k}(\mu, \nu)$$

and the result follows.

**Proposition 10.** Let  $\mathcal{X}$  and  $\mathcal{Y}$  be Polish spaces and  $\mathbf{c} := (c_i)_{1 \leq i \leq N}$  a family of nonnegative lower semi-continuous costs on  $\mathcal{X} \times \mathcal{Y}$ . We suppose that, for all  $i$ ,  $c_i = N \times c_1$ . Then for any  $(\mu, \nu) \in \mathcal{M}_+^1(\mathcal{X}) \times \mathcal{M}_+^1(\mathcal{Y})$

$$\text{EOT}_{\mathbf{c}}(\mu, \nu) = \text{EOT}_{c_1}(\mu, \nu) = W_{c_1}(\mu, \nu). \quad (27)$$

**Proof.** Let  $c := (c_i)_{1 \leq i \leq N}$  such that for all  $i$ ,  $c_i = c_1$ . for all  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  and  $\lambda \in \Delta_N^+$ , we have:

$$c_\lambda(x, y) := \min_i (\lambda_i c_i(x, y)) = \min_i (\lambda_i) c_1(x, y)$$

Therefore we obtain from Lemma 4 that

$$\text{EOT}_c(\mu, \nu) = \sup_{\lambda \in \Delta_N^+} W_{c_\lambda}(\mu, \nu) \quad (28)$$

But we also have that:

$$\begin{aligned} W_{c_\lambda}(\mu, \nu) &= \inf_{\gamma \in \Gamma(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} \min_i (\lambda_i c_i(x, y)) d\gamma(x, y) \\ &= \min_i (\lambda_i) \inf_{\gamma \in \Gamma(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} c_1(x, y) d\gamma(x, y) \\ &= \min_i (\lambda_i) W_{c_1}(\mu, \nu) \end{aligned}$$

Finally by taking the supremum over  $\lambda \in \Delta_N^+$  we conclude the proof.

### C.3 Regularized EOT tends to EOT

**Proposition 11.** For  $(\mu, \nu) \in \mathcal{M}_+^1(\mathcal{X}) \times \mathcal{M}_+^1(\mathcal{Y})$  we have  $\lim_{\varepsilon \rightarrow 0} \text{EOT}_{\mathbf{c}}^\varepsilon(\mu, \nu) = \text{EOT}_{\mathbf{c}}(\mu, \nu)$ .

**Proof.** Let  $(\varepsilon_l = (\varepsilon_{l,1}, \dots, \varepsilon_{l,N}))_l$  a sequence converging to 0. Let  $\gamma_l = (\gamma_{l,1}, \dots, \gamma_{l,N})$  be the optimum of  $\text{EOT}_{\mathbf{c}}^{\varepsilon_l}(\mu, \nu)$ . By Lemma 1, up to an extraction,  $\gamma_l \rightarrow \gamma^* = (\gamma_1^*, \dots, \gamma_N^*) \in \Gamma_{\mu, \nu}^N$ . Let now  $\gamma = (\gamma_1, \dots, \gamma_N)$  be the optimum of  $\text{EOT}_{\mathbf{c}}(\mu, \nu)$ . By optimality of  $\gamma$  and  $\gamma_l$ , for all  $i$ :

$$0 \leq \int c_i d\gamma_{l,i} - \int c_i d\gamma_i \leq \sum_i \varepsilon_{l,i} (\text{KL}(\gamma_i \| \mu \otimes \nu) - \text{KL}(\gamma_{l,i} \| \mu \otimes \nu))$$

By lower semi continuity of  $\text{KL}(\cdot \| \mu \otimes \nu)$  and by taking the limit inferior as  $l \rightarrow \infty$ , we get for all  $i$ ,  $\liminf_{l \rightarrow \infty} \int c_i d\gamma_{l,i} = \int c_i d\gamma_i$ . Moreover by continuity of  $\gamma \rightarrow \int c_i d\gamma_i$  we therefore obtain that for all  $i$ ,  $\int c_i d\gamma_i^* \leq \int c_i d\gamma_i$ . Then by optimality of  $\gamma$  the result follows.

### C.4 Projected Accelerated Gradient Descent

**Proposition 12.** Let  $a \in \Delta_N^+$  and  $b \in \Delta_m^+$  and  $\mathbf{C} := (C_i)_{1 \leq i \leq N} \in (\mathbb{R}^{n \times m})^N$  be  $N$  cost matrices and  $\varepsilon := (\varepsilon, \dots, \varepsilon)$  where  $\varepsilon > 0$ . Then by denoting  $K_i^{\lambda_i} = \exp(-\lambda_i C_i / \varepsilon)$ , we have

$$\widehat{\text{EOT}}_{\mathbf{C}}^\varepsilon(a, b) = \sup_{\lambda \in \Delta_N^+} \sup_{f \in \mathbb{R}^n, g \in \mathbb{R}^m} F_{\mathbf{C}}^\varepsilon(\lambda, f, g) := \langle f, a \rangle + \langle g, b \rangle - \varepsilon \left[ \log \left( \sum_{i=1}^N \langle e^{f/\varepsilon}, K_i^{\lambda_i} e^{g/\varepsilon} \rangle \right) + 1 \right].$$

Moreover,  $F_{\mathbf{C}}^\varepsilon$  is concave, differentiable and  $\nabla F$  is  $\frac{\max \left( \max_{1 \leq i \leq N} \|C_i\|_\infty^2, 2N \right)}{\varepsilon}$  Lipschitz-continuous on  $\mathbb{R}^N \times \mathbb{R}^n \times \mathbb{R}^m$ .

**Proof.** Let  $\mathcal{Q} := \left\{ P := (P_1, \dots, P_N) \in (\mathbb{R}_+^{n \times m})^N : \sum_{k=1}^N \sum_{i,j} P_k^{i,j} = 1 \right\}$ . Note that  $\Gamma_{a,b}^N \subset \mathcal{Q}$ , therefore from the

primal formulation of the problem we have that

$$\begin{aligned}\widehat{\text{EOT}}_{\mathbf{C}}^{\varepsilon}(a, b) &= \sup_{\lambda \in \Delta_N^+} \inf_{P \in \Gamma_{a,b}^N} \sum_{i=1}^N \lambda_i \langle P_i, C_i \rangle - \varepsilon H(P_i) \\ &= \sup_{\lambda \in \Delta_N^+} \inf_{P \in \mathcal{Q}} \sup_{f \in \mathbb{R}^n, g \in \mathbb{R}^m} \sum_{i=1}^N \lambda_i \langle P_i, C_i \rangle - \varepsilon H(P_i) \\ &\quad + f^T \left( a - \sum_i P_i \mathbf{1}_m \right) + g^T \left( b - \sum_i P_i^T \mathbf{1}_n \right)\end{aligned}$$

The constraints are qualified for this convex problem, hence by Slater's sufficient condition (Boyd et al., 2004, Section 5.2.3), strong duality holds. Therefore we have

$$\begin{aligned}\widehat{\text{EOT}}_{\mathbf{C}}^{\varepsilon}(a, b) &= \sup_{\lambda \in \Delta_N^+} \sup_{f \in \mathbb{R}^n, g \in \mathbb{R}^m} \inf_{P \in \mathcal{Q}} \sum_{i=1}^N \lambda_i \langle P_i, C_i \rangle - \varepsilon H(P_i) \\ &\quad + f^T \left( a - \sum_i P_i \mathbf{1}_m \right) + g^T \left( b - \sum_i P_i^T \mathbf{1}_n \right) \\ &= \sup_{\lambda \in \Delta_N^+} \sup_{f \in \mathbb{R}^n, g \in \mathbb{R}^m} \langle f, a \rangle + \langle g, b \rangle \\ &\quad + \inf_{P \in \mathcal{Q}} \sum_{k=1}^N \sum_{i,j} P_k^{i,j} \left( \lambda_k C_k^{i,j} + \varepsilon \left( \log(P_k^{i,j}) - 1 \right) - f_i - g_j \right)\end{aligned}$$

Let us now focus on the following problem:

$$\inf_{P \in \mathcal{Q}} \sum_{k=1}^N \sum_{i,j} P_k^{i,j} \left( \lambda_k C_k^{i,j} + \varepsilon \left( \log(P_k^{i,j}) - 1 \right) - f_i - g_j \right)$$

Note that for all  $i, j, k$  and some small  $\delta$ ,

$$P_k^{i,j} \left( \lambda_k C_k^{i,j} - \varepsilon \left( \log(P_k^{i,j}) - 1 \right) - f_i - g_j \right) < 0$$

if  $P_k^{i,j} \in (0, \delta)$  and this quantity goes to 0 as  $P_k^{i,j}$  goes to 0. Therefore  $P_k^{i,j} > 0$  and the problem becomes

$$\inf_{P > 0} \sup_{\nu \in \mathbb{R}} \sum_{k=1}^N \sum_{i,j} P_k^{i,j} \left( \lambda_k C_k^{i,j} + \varepsilon \left( \log(P_k^{i,j}) - 1 \right) - f_i - g_j \right) + \nu \left( \sum_{k=1}^N \sum_{i,j} P_k^{i,j} - 1 \right).$$

The solution to this problem is for all  $k \in \{1, \dots, N\}$ ,

$$P_k = \frac{\exp \left( \frac{f \mathbf{1}_n^T + \mathbf{1}_m g^T - \lambda_k C_k}{\varepsilon} \right)}{\sum_{k=1}^N \sum_{i,j} \exp \left( \frac{f_i + g_j - \lambda_k C_k^{i,j}}{\varepsilon} \right)}$$

Therefore we obtain that

$$\begin{aligned}\widehat{\text{EOT}}_{\mathbf{C}}^{\varepsilon}(a, b) &= \sup_{\lambda \in \Delta_N^+} \sup_{f \in \mathbb{R}^n, g \in \mathbb{R}^m} \langle f, a \rangle + \langle g, b \rangle \\ &\quad - \varepsilon \sum_{k=1}^N \sum_{i,j} P_k^{i,j} \left[ \log \left( \sum_{k=1}^N \sum_{i,j} \exp \left( \frac{f_i + g_j - \lambda_k C_k^{i,j}}{\varepsilon} \right) \right) + 1 \right] \\ &= \sup_{\lambda \in \Delta_N^+} \sup_{f \in \mathbb{R}^n, g \in \mathbb{R}^m} \langle f, a \rangle + \langle g, b \rangle - \varepsilon \left[ \log \left( \sum_{k=1}^N \sum_{i,j} \exp \left( \frac{f_i + g_j - \lambda_k C_k^{i,j}}{\varepsilon} \right) \right) + 1 \right].\end{aligned}$$

From now on, we denote for all  $\lambda \in \Delta_N^+$

$$\widehat{\text{EOT}}_{\mathbf{C}}^{\varepsilon, \lambda}(a, b) := \inf_{P \in \Gamma_{a, b}^N} \sum_{i=1}^N \lambda_i \langle P_i, C_i \rangle - \varepsilon H(P_i)$$

$$\widehat{\text{EOT}}_{\mathbf{C}}^{\varepsilon, \lambda}(a, b) := \sup_{f \in \mathbb{R}^n, g \in \mathbb{R}^m} \langle f, a \rangle + \langle g, b \rangle - \varepsilon \left[ \log \left( \sum_{k=1}^N \sum_{i, j} \exp \left( \frac{f_i + g_j - \lambda_k C_k^{i, j}}{\varepsilon} \right) \right) + 1 \right]$$

which has just been shown to be dual and equal. Thanks to (Nesterov, 2005, Theorem 1), as for all  $\lambda \in \mathbb{R}^N$ ,  $P \in \Gamma_{a, b}^N \rightarrow \sum_{i=1}^N \lambda_i \langle P_i, C_i \rangle - \varepsilon H(P_i)$  is  $\varepsilon$ -strongly convex, then for all  $\lambda \in \mathbb{R}^N$ ,  $(f, g) \rightarrow \nabla_{(f, g)} F(\lambda, f, g)$  is  $\frac{\|A\|_{1 \rightarrow 2}^2}{\varepsilon}$  Lipschitz-continuous where  $A$  is the linear operator of the equality constraints of the primal problem. Moreover this norm is equal to the maximum Euclidean norm of a column of  $A$ . By definition, each column of  $A$  contains only  $2N$  non-zero elements, which are equal to one. Hence,  $\|A\|_{1 \rightarrow 2} = \sqrt{2N}$ . Let us now show that for all  $(f, g) \in \mathbb{R}^n \times \mathbb{R}^m$   $\lambda \in \mathbb{R}^N \rightarrow \nabla_{\lambda} F(\lambda, f, g)$  is also Lipschitz-continuous. Indeed we remarks that

$$\frac{\partial^2 F}{\partial \lambda_q \partial \lambda_k} = \frac{1}{\varepsilon \nu^2} [\sigma_{q,1}(\lambda) \sigma_{k,1}(\lambda) - \nu(\sigma_{k,2}(\lambda) \mathbb{1}_{k=q})]$$

where  $\mathbb{1}_{k=q} = 1$  iff  $k = q$  and 0 otherwise, for all  $k \in \{1, \dots, N\}$  and  $p \geq 1$

$$\sigma_{k,p}(\lambda) = \sum_{i,j} (C_k^{i,j})^p \exp \left( \frac{f_i + g_j - \lambda_k C_k^{i,j}}{\varepsilon} \right)$$

$$\nu = \sum_{k=1}^N \sum_{i,j} \exp \left( \frac{f_i + g_j - \lambda_k C_k^{i,j}}{\varepsilon} \right).$$

Let  $v \in \mathbb{R}^N$ , and by denoting  $\nabla_{\lambda}^2 F$  the Hessian of  $F$  with respect to  $\lambda$  for fixed  $f, g$  we obtain first that

$$\begin{aligned} v^T \nabla_{\lambda}^2 F v &= \frac{1}{\varepsilon \nu^2} \left[ \left( \sum_{k=1}^N v_k \sigma_{q,1}(\lambda) \right)^2 - \nu \sum_{k=1}^N v_k^2 \sigma_{k,2} \right] \\ &\leq \frac{1}{\varepsilon \nu^2} \left( \sum_{k=1}^N v_k \sigma_{q,1}(\lambda) \right)^2 \\ &\quad - \frac{1}{\varepsilon \nu^2} \left( \sum_{k=1}^N |v_k| \sqrt{\sum_{i,j} \exp \left( \frac{f_i + g_j - \lambda_k C_k^{i,j}}{\varepsilon} \right)} \sqrt{\sum_{i,j} (C_k^{i,j})^2 \exp \left( \frac{f_i + g_j - \lambda_k C_k^{i,j}}{\varepsilon} \right)} \right)^2 \\ &\leq \frac{1}{\varepsilon \nu^2} \left[ \left( \sum_{k=1}^N v_k \sigma_{q,1}(\lambda) \right)^2 - \left( \sum_{k=1}^N |v_k| \sum_{i,j} |C_k^{i,j}| \exp \left( \frac{f_i + g_j - \lambda_k C_k^{i,j}}{\varepsilon} \right) \right)^2 \right] \\ &\leq 0 \end{aligned}$$

Indeed the last two inequalities come from Cauchy Schwartz. Moreover we have

$$\begin{aligned} \frac{1}{\varepsilon \nu^2} \left[ \left( \sum_{k=1}^N v_k \sigma_{q,1}(\lambda) \right)^2 - \nu \sum_{k=1}^N v_k^2 \sigma_{k,2} \right] &= v^T \nabla_{\lambda}^2 F v \leq 0 \\ &\quad - \frac{\sum_{k=1}^N v_k^2 \sigma_{k,2}}{\varepsilon \nu} \leq \\ &\quad - \frac{\sum_{k=1}^N v_k^2 \max_{1 \leq i \leq N} (\|C_i\|_{\infty}^2)}{\varepsilon} \leq \end{aligned}$$

Therefore we deduce that  $\lambda \in \mathbb{R}^N \rightarrow \nabla_\lambda F(\lambda, f, g)$  is  $\frac{\max_{1 \leq i \leq N} (\|C_i\|_\infty^2)}{\varepsilon}$  Lipschitz-continuous, hence  $\nabla F(\lambda, f, g)$  is  $\frac{\max_{1 \leq i \leq N} (\|C_i\|_\infty^2, 2N)}{\varepsilon}$  Lipschitz-continuous on  $\mathbb{R}^N \times \mathbb{R}^n \times \mathbb{R}^m$ .

Denote  $L := \frac{\max_{1 \leq i \leq N} (\|C_i\|_\infty^2, 2N)}{\varepsilon}$  the Lipschitz constant of  $F_{\mathbf{C}}^\varepsilon$ . Moreover for all  $\lambda \in \mathbb{R}^N$ , let  $\text{Proj}_{\Delta_N^+}(\lambda)$  the unique solution of the following optimization problem

$$\min_{x \in \Delta_N^+} \|x - \lambda\|_2^2. \quad (29)$$

Let us now introduce the following algorithm.

---

**Algorithm 2** Accelerated Projected Gradient Ascent Algorithm
 

---

**Input:**  $\mathbf{C} = (C_i)_{1 \leq i \leq N}$ ,  $a$ ,  $b$ ,  $\varepsilon$ ,  $L$

**Init:**  $f^{-1} = f^0 \leftarrow \mathbf{0}_n$ ;  $g^{-1} = g^0 \leftarrow \mathbf{0}_m$ ;  $\lambda^{-1} = \lambda^0 \leftarrow (1/N, \dots, 1/N) \in \mathbb{R}^N$

**for**  $k = 1, 2, \dots$  **do**

$$\begin{aligned} & (v, w, z)^T \leftarrow (\lambda^{k-1}, f^{k-1}, g^{k-1})^T + \frac{k-2}{k+1} ((\lambda^{k-1}, f^{k-1}, g^{k-1})^T - (\lambda^{k-2}, f^{k-2}, g^{k-2})^T); \\ & \lambda^k \leftarrow \text{Proj}_{\Delta_N^+} \left( v + \frac{1}{L} \nabla_\lambda F_{\mathbf{C}}^\varepsilon(v, w, z) \right); \\ & (g^k, f^k)^T \leftarrow (w, z)^T + \frac{1}{L} \nabla_{(f,g)} F_{\mathbf{C}}^\varepsilon(v, w, z). \end{aligned}$$

**end**

**Result:**  $\lambda, f, g$

---

Beck and Teboulle (2009); Tseng (2008) give us that the accelerated projected gradient ascent algorithm achieves the optimal rate for first order methods of  $\mathcal{O}(1/k^2)$  for smooth functions. To perform the projection we use the algorithm proposed in Shalev-Shwartz and Singer (2006) which finds the solution of (29) after  $\mathcal{O}(N \log(N))$  algebraic operations (Wang and Carreira-Perpinan, 2013).

### C.5 Fair cutting cake problem

Let  $\mathcal{X}$ , be a set representing a cake. The aim of the cutting cake problem is to divide it in  $\mathcal{X}_1, \dots, \mathcal{X}_N$  disjoint sets among the  $N$  individuals. The utility for a single individual  $i$  for a slice  $S$  is denoted  $V_i(S)$ . It is often assumed that  $V_i(\mathcal{X}) = 1$  and that  $V_i$  is additive for disjoint sets. There exists many criteria to assess fairness for a partition  $\mathcal{X}_1, \dots, \mathcal{X}_N$  such as proportionality ( $V_i(\mathcal{X}_i) \geq 1/N$ ), envy-freeness ( $V_i(\mathcal{X}_i) \geq V_i(\mathcal{X}_j)$ ) or equitability ( $V_i(\mathcal{X}_i) = V_j(\mathcal{X}_j)$ ). A possible problem to solve equitability and proportionality in the cutting cake problem is the following:

$$\inf_{\substack{\mathcal{X}_1, \dots, \mathcal{X}_N \\ \sqcup_{i=1}^N \mathcal{X}_i = \mathcal{X}}} \max_i V_i(\mathcal{X}_i) \quad (30)$$

Note that here we do not want to solve the problem under equality constraints since the problem might not be well defined. Moreover the existence of the optimum is not immediate. A natural relaxation of this problem is when there is a divisible quantity of each element of the cake ( $x \in \mathcal{X}$ ). In that case, the cake is no more a set but rather a distribution on this set  $\mu$ . Following the primal formulation of EOT, it is clear that it is a relaxation of the cutting cake problem where the goal is to divide the cake viewed as a distribution. For the cutting cake problem with two cakes  $\mathcal{X}$  and  $\mathcal{Y}$ , the problem can be cast as follows:

$$\inf_{\substack{\mathcal{X}_1, \dots, \mathcal{X}_N \text{ s.t. } \sqcup_{i=1}^N \mathcal{X}_i = \mathcal{X} \\ \mathcal{Y}_1, \dots, \mathcal{Y}_N \text{ s.t. } \sqcup_{i=1}^N \mathcal{Y}_i = \mathcal{Y}}} \max_i V_i(\mathcal{X}_i, \mathcal{Y}_i) \quad (31)$$

Here EOT is the relaxation of this problem where we split the cakes viewed as distributions instead of sets themselves. Note that in this problem, the utility of the agents are coupled.

## D Illustrations and Experiments

### D.1 Primal Formulation

Here we show the couplings obtained when we consider three negative costs  $\tilde{c}_i$  which corresponds to the situation where we aim to obtain a fair division of goods between three agents. Moreover we show the couplings obtained according to the transport viewpoint where we consider the opposite of these three negative cost functions, i.e.  $c_i := -\tilde{c}_i$ . We can see that the couplings obtained in the two situations are completely different, which is expected. Indeed in the fair division problem, we aim at finding couplings which maximize the total utility of each agent ( $\int c_i d\gamma_i^1$ ) while ensuring that their are equal while in the other case, we aim at finding couplings which minimize the total transportation cost of each agent ( $\int c_i d\gamma_i^2$ ) while ensuring that their are equal. Obviously we always have that

$$\forall i \quad \int c_i d\gamma_i^2 \leq \int c_i d\gamma_i^1.$$

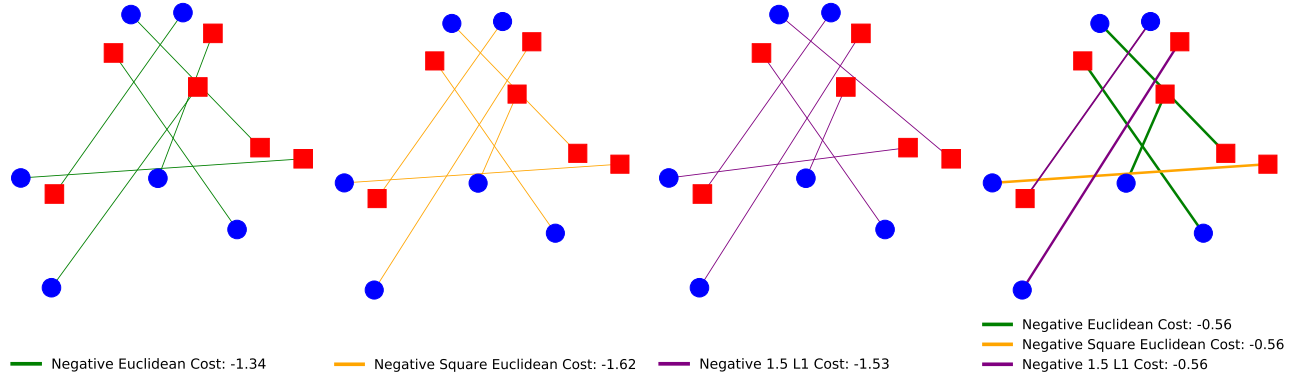


Figure 4: Comparison of the optimal couplings obtained from standard OT for three different costs and EOT in case of negative costs (i.e. utilities). Blue dots and red squares represent the locations of two discrete uniform measures. *Left, middle left, middle right*: Kantorovich couplings between the two measures for negative Euclidean cost ( $-\|\cdot\|_2$ ), negative square Euclidean cost ( $-\|\cdot\|_2^2$ ) and negative 1.5 L1 norm ( $-\|\cdot\|_1^{1.5}$ ) respectively. *Right*: Equitable and optimal division of the resources between the  $N = 3$  different negative costs (i.e. utilities) given by EOT. Note that the partition between the agents is equitable (i.e. utilities are equal) and proportional (i.e. utilities are larger than  $1/N$ ).

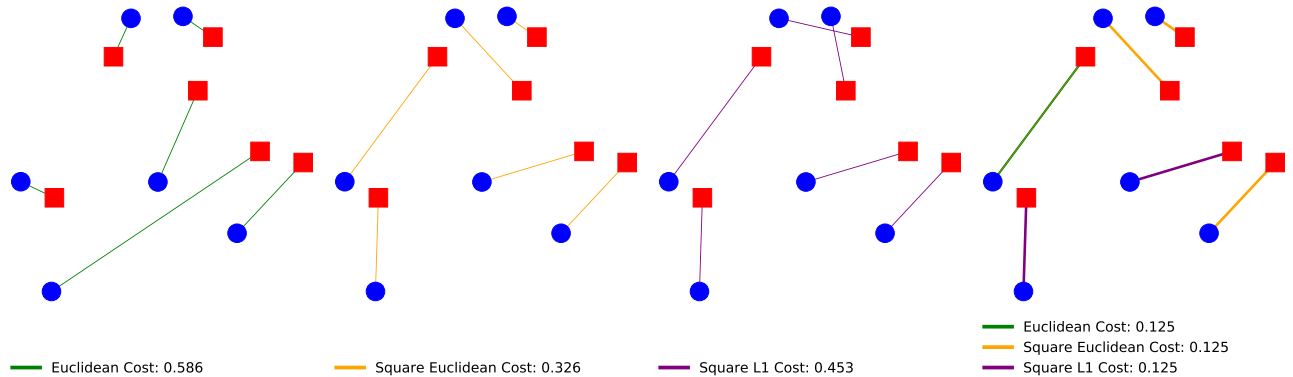


Figure 5: Comparison of the optimal couplings obtained from standard OT for three different costs and EOT in case of postive costs. Blue dots and red squares represent the locations of two discrete uniform measures. *Left, middle left, middle right*: Kantorovich couplings between the two measures for Euclidean cost ( $\|\cdot\|_2$ ), square Euclidean cost ( $\|\cdot\|_2^2$ ) and 1.5 L1 norm ( $\|\cdot\|_1^{1.5}$ ) respectively. *Right*: transport couplings of EOT solving Eq. (1). Note that each cost contributes equally and its contribution is lower than the smallest OT cost.



## D.2 Dual Formulation

Here we show the dual variables obtained in the exact same settings as in the primal illustrations. Figure 6 shows the dual associated to the primal problem exposed in Figure 4 and Figure 7 shows the dual associated to the primal problem exposed in Figure 5.

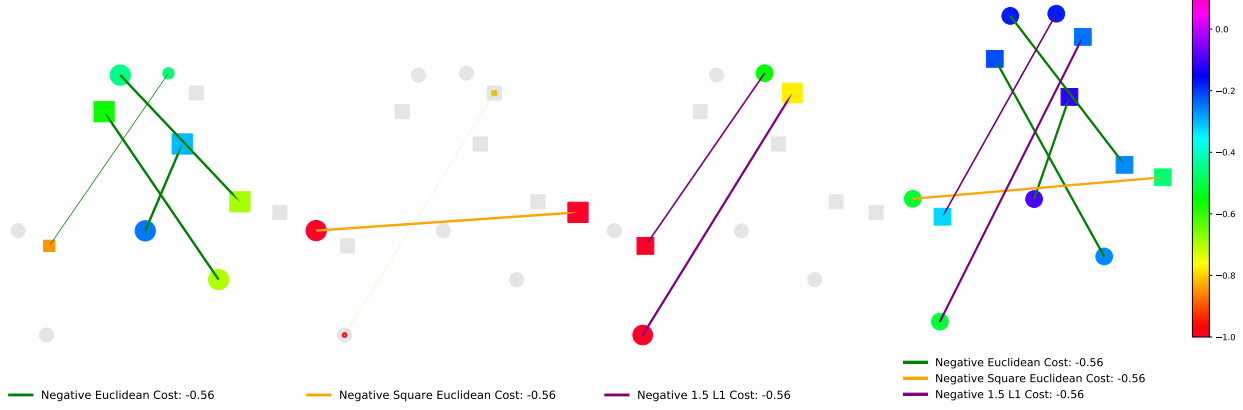


Figure 6: *Left, middle left, middle right*: the size of dots and squares is proportional to the weight of their representing atom in the distributions  $\mu_k^*$  and  $\nu_k^*$  respectively. The utilities  $f_k^*$  and  $g_k^*$  for each point in respectively  $\mu_k^*$  and  $\nu_k^*$  are represented by the color of dots and squares according to the color scale on the right hand side. The gray dots and squares correspond to the points that are ignored by agent  $k$  in the sense that there is no mass or almost no mass in distributions  $\mu_k^*$  or  $\nu_k^*$ . *Right*: the size of dots and squares are uniform since they correspond to the weights of uniform distributions  $\mu$  and  $\nu$  respectively. The values of  $f^*$  and  $g^*$  are given also by the color at each point. Note that each agent gets exactly the same total utility, corresponding exactly to EOT. This value can be computed using dual formulation (5) and for each figure it equals the sum of the values (encoded with colors) multiplied by the weight of each point (encoded with sizes).

**Transport viewpoint of the Dual Formulation.** Assume that the  $N$  agents are not able to solve the primal problem (1) which aims at finding the cheapest equitable partition of the work among the  $N$  agents for transporting the distributions of goods  $\mu$  to the distributions of stores  $\nu$ . Moreover assume that there is an external agent who can do the transportation work for them with the following pricing scheme: he or she splits the logistic task into that of collecting and then delivering the goods, and will apply a collection price  $\tilde{f}(x)$  for one unit of good located at  $x$  (no matter where that unit is sent to), and a delivery price  $\tilde{g}(y)$  for one unit to the location  $y$  (no matter from which place that unit comes from). Then the external agent for transporting some goods  $\mu$  to some stores  $\nu$  will charge  $\int_{x \in \mathcal{X}} \tilde{f}(x) d\mu(x) + \int_{y \in \mathcal{Y}} \tilde{g}(y) d\nu(y)$ . However he or she has the constraint that the pricing must be equitable among the agents and therefore wants to ensure that each agent will pay exactly  $\frac{1}{N} \int_{x \in \mathcal{X}} \tilde{f}(x) d\mu(x) + \int_{y \in \mathcal{Y}} \tilde{g}(y) d\nu(y)$ . Denote  $f = \frac{\tilde{f}}{N}$ ,  $g = \frac{\tilde{g}}{N}$  and therefore the price paid by each agent becomes  $\int_{x \in \mathcal{X}} f(x) d\mu(x) + \int_{y \in \mathcal{Y}} g(y) d\nu(y)$ . Moreover, to ensure that each agent will not pay more than he would if he was doing the job himself or herself, he or she must guarantee that for all  $\lambda \in \Delta_N^+$ , the pricing scheme  $(f, g)$  satisfies:

$$f \oplus g \leq \min(\lambda_i c_i).$$

Indeed under this constraint, it is easy for the agents to check that they will never pay more than what they would pay if they were doing the transportation task as we have

$$\int_{x \in \mathcal{X}} f(x) d\mu(x) + \int_{y \in \mathcal{Y}} g(y) d\nu(y) \leq \int_{\mathcal{X} \times \mathcal{Y}} \min_i(\lambda_i c_i) d\gamma$$

which holds for every  $\gamma$  in particular for  $\gamma^* = \sum_{i=1}^N \gamma_i^*$  optimal solution of the primal problem (1) from which follows

$$\begin{aligned} \int_{x \in \mathcal{X}} f(x) d\mu(x) + \int_{y \in \mathcal{Y}} g(y) d\nu(y) &\leq \sum_{i=1}^N \int_{\mathcal{X} \times \mathcal{Y}} \min_i(\lambda_i c_i) d\gamma_i^* \\ &\leq \sum_{i=1}^N \lambda_i \int_{\mathcal{X} \times \mathcal{Y}} c_i d\gamma_i^* \\ &= \text{EOT}_{\mathbf{c}}(\mu, \nu) \end{aligned}$$

Therefore the external agent aims to maximise his or her selling price under the above constraints which is exactly the dual formulation of our problem.

Another interpretation of the dual problem when the cost are non-negative can be expressed as follows. Let us introduce the subset of  $(\mathcal{C}^b(\mathcal{X}) \times \mathcal{C}^b(\mathcal{Y}))^N$ :

$$\mathcal{G}_{\mathbf{c}}^N := \{(f_k, g_k)_{k=1}^N \text{ s.t. } \forall k, f_k \oplus g_k \leq c_k\}$$

Let us now show the following reformulation of the problem. See Appendix D.2 for the proof.

**Proposition 13.** *Under the same assumptions of Proposition 1, we have*

$$\begin{aligned} \text{EOT}_{\mathbf{c}}(\mu, \nu) &= \sup_{(f_k, g_k)_{k=1}^N \in \mathcal{G}_{\mathbf{c}}^N} \inf_{\substack{t \in \mathbb{R} \\ (\mu_k, \nu_k)_{k=1}^N \in \Upsilon_{\mu, \nu}^N}} t \\ &\text{ s.t. } \forall k, \int f_k d\mu_k + \int g_k d\nu_k = t \end{aligned} \quad (32)$$

**Proof.** Let us first introduce the following Lemma which guarantees that compactity of  $\Upsilon_{\mu, \nu}^N$  for the weak topology.

**Lemma 6.** *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be Polish spaces, and  $\mu$  and  $\nu$  two probability measures respectively on  $\mathcal{X}$  and  $\mathcal{Y}$ . Then  $\Upsilon_{\mu, \nu}^N$  is sequentially compact for the weak topology induced by  $\|\gamma\| = \max_{i=1, \dots, N} \|\mu_i\|_{\text{TV}} + \|\nu_i\|_{\text{TV}}$ .*

**Proof.** Let  $(\gamma^n)_{n \geq 0}$  a sequence in  $\Upsilon_{\mu, \nu}^N$ , and let us denote for all  $n \geq 0$ ,  $\gamma^n = (\mu_i^n, \nu_i^n)_{i=1}^N$ . We first remarks that for all  $i \in \{1, \dots, N\}$  and  $n \geq 0$ ,  $\|\mu_i^n\|_{\text{TV}} \leq 1$  and  $\|\nu_i^n\|_{\text{TV}} \leq 1$  therefore for all  $i \in \{1, \dots, N\}$ ,  $(\mu_i^n)_{n \geq 0}$  and  $(\nu_i^n)_{n \geq 0}$  are uniformly bounded. Moreover as  $\{\mu\}$  and  $\{\nu\}$  are tight, for any  $\delta > 0$ , there exists  $K \subset \mathcal{X}$  and  $L \subset \mathcal{Y}$  compact such that  $\mu(K^c) \leq \delta$  and  $\nu(L^c) \leq \delta$ . Then, we obtain that for any for all  $i \in \{1, \dots, N\}$ ,  $\mu_i^n(K^c) \leq \delta$  and  $\nu_i^n(L^c) \leq \delta$ . Therefore, for all  $i \in \{1, \dots, N\}$ ,  $(\mu_i^n)_{n \geq 0}$  and  $(\nu_i^n)_{n \geq 0}$  are tight and uniformly bounded and Prokhorov's theorem (Dupuis and Ellis, 2011, Theorem A.3.15) guarantees for all  $i \in \{1, \dots, N\}$ ,  $(\mu_i^n)_{n \geq 0}$  and  $(\nu_i^n)_{n \geq 0}$  admit a weakly convergent subsequence. By extracting a common convergent subsequence, we obtain that  $(\gamma^n)_{n \geq 0}$  admits a weakly convergent subsequence. By continuity of the projection, the limit also lives in  $\Upsilon_{\mu, \nu}^N$  and the result follows.

We can now prove the Proposition. We have that for any  $\lambda \in \Delta_N$

$$\begin{aligned} &\sup_{(f, g) \in \mathcal{F}_{\mathbf{c}}^\lambda} \int_{x \in \mathcal{X}} f(x) d\mu(x) + \int_{y \in \mathcal{Y}} g(y) d\nu(y) \\ &\leq \sup_{(f_k, g_k)_{k=1}^N \in \mathcal{G}_{\mathbf{c}}^N} \inf_{(\mu_k, \nu_k)_{k=1}^N \in \Upsilon_{\mu, \nu}^N} \sum_{k=1}^N \lambda_k \left[ \int_{x \in \mathcal{X}} f_k(x) d\mu_k(x) + \int_{y \in \mathcal{Y}} g_k(y) d\nu_k(y) \right] \\ &\leq \text{EOT}_{\mathbf{c}}(\mu, \nu) \end{aligned}$$

Then by taking the supremum over  $\lambda \in \Delta_N$ , and by applying Theorem 1 we obtain that

$$\text{EOT}_{\mathbf{c}}(\mu, \nu) = \sup_{\lambda \in \Delta_N} \sup_{(f_k, g_k)_{k=1}^N \in \mathcal{G}_{\mathbf{c}}^N} \inf_{(\mu_k, \nu_k)_{k=1}^N \in \Upsilon_{\mu, \nu}^N} \sum_{k=1}^N \lambda_k \left[ \int_{x \in \mathcal{X}} f_k(x) d\mu_k(x) + \int_{y \in \mathcal{Y}} g_k(y) d\nu_k(y) \right]$$

Let  $\mathcal{G}_{\mathbf{c}}^N$  and  $\Upsilon_{\mu, \nu}^N$  be endowed respectively with the uniform norm and the norm defined in Lemma 6. Note that the objective is linear and continuous with respect to  $(\mu_k, \nu_k)_{k=1}^N$  and also  $(f_k, g_k)_{k=1}^N$ . Moreover the spaces  $\mathcal{G}_{\mathbf{c}}^N$

and  $\Upsilon_{\mu,\nu}^N$  are clearly convex. Finally thanks to Lemma 6,  $\Upsilon_{\mu,\nu}^N$  is compact with respect to the weak topology we can apply Sion's theorem Sion (1958) and we obtain that

$$\text{EOT}_{\mathbf{c}}(\mu, \nu) = \sup_{(f_k, g_k)_{k=1}^N \in \mathcal{G}_{\mathbf{c}}^N} \inf_{(\mu_k, \nu_k)_{k=1}^N \in \Upsilon_{\mu,\nu}^N} \sup_{\lambda \in \Delta_N} \sum_{k=1}^N \lambda_k \left[ \int_{x \in \mathcal{X}} f_k(x) d\mu_k(x) + \int_{y \in \mathcal{Y}} g_k(y) d\nu_k(y) \right]$$

Let us now fix  $(f_k, g_k)_{k=1}^N \in \mathcal{G}_{\mathbf{c}}^N$  and  $(\mu_k, \nu_k)_{k=1}^N \in \Upsilon_{\mu,\nu}^N$ , therefore we have:

$$\begin{aligned} & \sup_{\lambda \in \Delta_N} \sum_{k=1}^N \lambda_k \left[ \int_{x \in \mathcal{X}} f_k(x) d\mu_k(x) + \int_{y \in \mathcal{Y}} g_k(y) d\nu_k(y) \right] \\ &= \sup_{\lambda} \inf_t t \times \left( 1 - \sum_{i=1}^N \lambda_i \right) + \sum_{k=1}^N \lambda_k \left[ \int_{x \in \mathcal{X}} f_k(x) d\mu_k(x) + \int_{y \in \mathcal{Y}} g_k(y) d\nu_k(y) \right] \\ &= \inf_t \sup_{\lambda} t + \sum_{k=1}^N \lambda_k \left[ \int_{x \in \mathcal{X}} f_k(x) d\mu_k(x) + \int_{y \in \mathcal{Y}} g_k(y) d\nu_k(y) - t \right] \\ &= \inf_t \left\{ t \text{ s.t. } \forall k, \int f_k d\mu_k + \int g_k d\nu_k = t \right\} \end{aligned}$$

where the inversion is possible as the Slater's conditions are satisfied and the result follows.

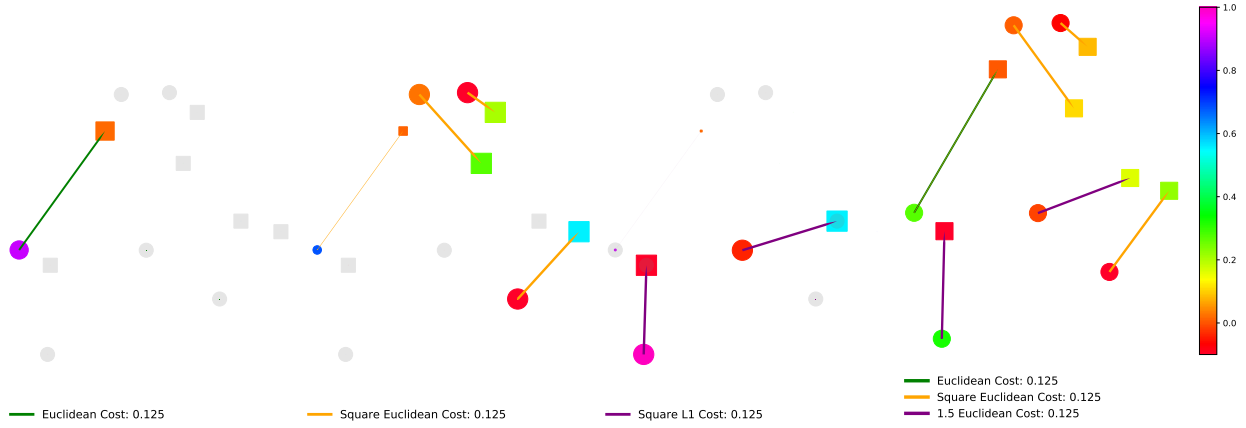


Figure 7: *Left, middle left, middle right:* the size of dots and squares is proportional to the weight of their representing atom in the distributions  $\mu_k^*$  and  $\nu_k^*$  respectively. The collection “cost”  $f_k^*$  for each point in  $\mu_k^*$ , and its delivery counterpart  $g_k^*$  in  $\nu_k^*$  are represented by the color of dots and squares according to the color scale on the right hand side. The gray dots and squares correspond to the points that are ignored by agent  $k$  in the sense that there is no mass or almost no mass in distributions  $\mu_k^*$  or  $\nu_k^*$ . *Right:* the size of dots and squares are uniform since they corresponds to the weights of uniform distributions  $\mu$  and  $\nu$  respectively. The values of  $f^*$  and  $g^*$  are given also by the color at each point. Note that each agent earns exactly the same amount of money, corresponding exactly EOT cost. This value can be computed using dual formulation (5) or its reformulation (32) and for each figure it equals the sum of the values (encoded with colors) multiplied by the weight of each point (encoded with sizes).

### D.3 Approximation of the Dudley Metric

Figure 8 illustrates the convergence of the entropic regularization approximation when  $\epsilon \rightarrow 0$ . To do so we plot the relative error from the ground truth defined as  $RE := \frac{EOT_\epsilon^\epsilon - \beta_d}{\beta_d}$  for different regularizations where  $\beta_d$  is obtained by solving the exact linear program and  $EOT_\epsilon^\epsilon$  is obtained by our proposed Alg. 1.

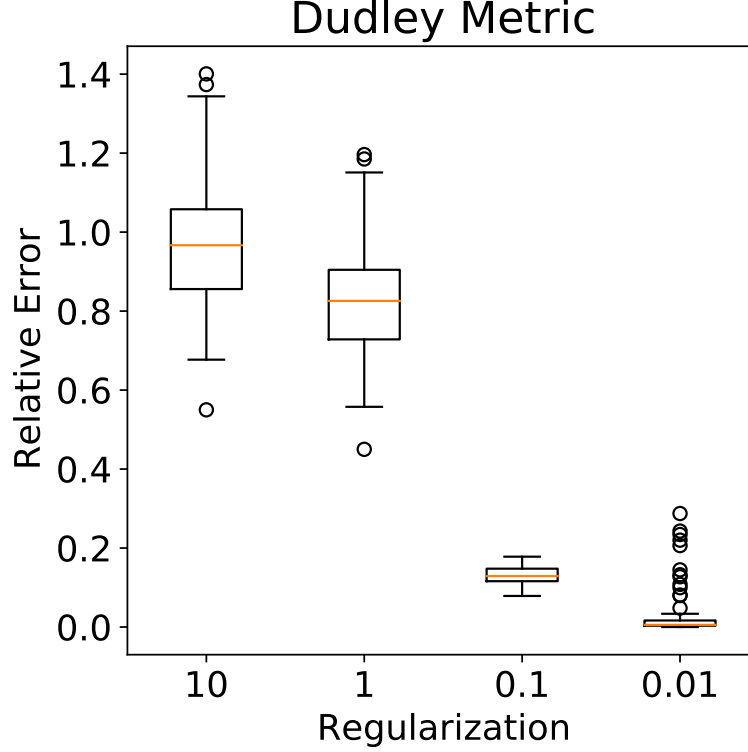


Figure 8: In this experiment, we draw 100 samples from two normal distributions and we plot the relative error from ground truth for different regularizations. We consider the case where two costs are involved:  $c_1 = 2 \times \mathbf{1}_{x \neq y}$ , and  $c_2 = d$  where  $d$  is the Euclidean distance. This case corresponds exactly to the Dudley metric (see Proposition 4). We remark that as  $\epsilon \rightarrow 0$ , the approximation error goes also to 0.