

Optimal Tensor Transport

Tanguy Kerdoncuff¹, Rémi Emonet¹, Michaël Perrot², Marc Sebban¹

¹ Univ Lyon, UJM-Saint-Etienne, CNRS, Institut d'Optique Graduate School, Laboratoire Hubert Curien UMR 5516, F-42023, Saint-Etienne, France

² Univ. Lille, Inria, CNRS, Centrale Lille, UMR 9189 - CRIStAL, F-59000 Lille, France
{tanguy.kerdoncuff, remi.emonet, marc.sebban}@univ-st-etienne.fr, michael.perrot@inria.fr

Abstract

Optimal Transport (OT) has become a popular tool in machine learning to align finite datasets typically lying in the same vector space. To expand the range of possible applications, Co-Optimal Transport (Co-OT) jointly estimates two distinct transport plans, one for the rows (points) and one for the columns (features), to match two data matrices that might use different features. On the other hand, Gromov Wasserstein (GW) looks for a single transport plan from two pairwise intra-domain distance matrices. Both Co-OT and GW can be seen as specific extensions of OT to more complex data. In this paper, we propose a unified framework, called Optimal Tensor Transport (OTT), which takes the form of a generic formulation that encompasses OT, GW and Co-OT and can handle tensors of any order by learning possibly multiple transport plans. We derive theoretical results for the resulting new distance and present an efficient way for computing it. We further illustrate the interest of such a formulation in Domain Adaptation and Comparison-based Clustering.

Introduction

Comparing two probability measures in the form of empirical distributions is at the core of many machine learning tasks. Optimal Transport (OT) (Villani 2008; Peyré, Cuturi et al. 2019) is a popular tool that allows such comparisons for datasets typically lying in a common vector space. Given two point clouds and a metric allowing to evaluate the transportation cost between two samples, the goal of OT is to learn the transport plan that minimizes the alignment cost between the two sets, resulting in the so-called Wasserstein distance. OT has been shown to be of great interest when dealing with machine learning tasks. For example, unsupervised Domain Adaptation (DA) aims at benefiting from labeled data of a source domain to classify examples drawn from a different but related target domain. The DA theory prompts us to reduce the shift between the source and the target distributions, a task that can be addressed by aligning the two datasets using OT (Courty et al. 2016, 2017; Shen et al. 2018; Damodaran et al. 2018). OT has also been successfully used in generative adversarial networks (GAN) (Goodfellow et al. 2014)

to minimize the divergence between training data and samples drawn from a generative model, leading to the WGAN (Arjovsky, Chintala, and Bottou 2017; Gulrajani et al. 2017).

In order to expand the range of possible applications, different variants have been proposed in the OT literature to tackle more complex settings. While the standard OT scenario assumes that the two datasets lie in the same feature space, Gromov Wasserstein (GW) (Memoli 2007; Mémoli 2011; Peyré, Cuturi, and Solomon 2016) extends the framework to incomparable spaces by allowing the alignment of two distributions when only the within-dataset pairwise distances are available. This approach is particularly well suited to deal with graphs described by their adjacency matrices (Xu et al. 2019; Xu, Luo, and Carin 2019; Chowdhury and Mémoli 2019). The GW discrepancy has been used efficiently in various applications such as heterogeneous DA (Yan et al. 2018), word translation (Alvarez-Melis and Jaakkola 2018) or GAN (Vayer et al. 2019; Bunne et al. 2019). Recently, Co-Optimal Transport (Co-OT) (Redko et al. 2020) extended the OT theory to datasets lying in different vector spaces. The idea is to jointly learn two transport plans. The first one aligns the examples as in standard OT while the second one aligns the most similar features. This has been shown to be of particular interest in heterogeneous DA and co-clustering.

Motivation and Contribution. While GW and Co-OT already cover a wide range of problems, we claim that many other scenarios are not covered by these two extensions. Let us suppose that both the source and target distributions are represented by a collection of graphs of the same size (in terms of nodes) but of different structure (in terms of edges). This is typically the case of two graphs evolving over time. In this case, the goal of OT would be to jointly align both the two collections of graphs and the nodes. It turns out that GW would be only able to handle the special case where there exist a known one-to-one correspondence between the graphs of the two collections. Another application is inspired from comparison-based learning. Let us consider a source and a target distribution represented by a set of users who watched movies, users providing a list of triplet comparisons of the form “movie x_i is closer to x_j than to x_k ”. In this case, neither GW nor Co-OT is able to align the two distributions because of the nature of this triplet-based representation. A last example comes from computer vision, where one may want to align two collections of images while preserving

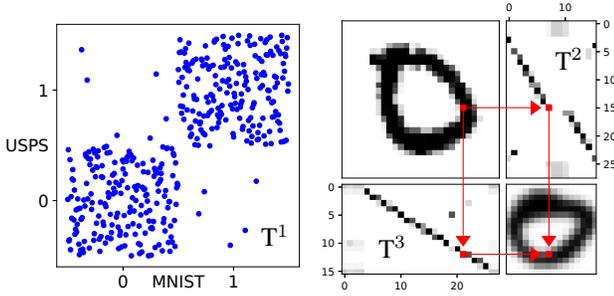


Figure 1: (Left) Transport plan T^1 between 400 images (only digits 0 and 1) of MNIST and USPS datasets; (Right) (top left) An example from MNIST and (bottom right) an example from USPS with a 90° right rotation; (top right) the OT plan T^2 between the rows of MNIST and USPS; (bottom left) the OT plan T^3 between the columns of MNIST and USPS; the arrows explain how to match the pixels between the two datasets using T^2 and T^3 obtained with OTT.

some inner structural information in rows and columns. It is worth noting that these three applications share a common characteristic: they can be represented in the form of third-order tensors. To solve OT tasks on such complex structures, it is necessary to design a framework generalizing the OT theory. This is the main contribution of this paper.

We propose *Optimal Tensor Transport* (OTT), a new OT formulation that can handle datasets represented as tensors of any order by potentially learning multiple transport plans. The underlying idea is to jointly match the different dimensions of each tensor with respect to their weights. Figure 1 illustrates this on a transportation problem between images from the MNIST (LeCun et al. 1998) to the USPS dataset (Friedman et al. 2001). Three transport plans are optimized in this scenario. T^1 is used to match the points (on the left), T^2 and T^3 preserve the structure by respectively mapping the pixel rows and pixel columns jointly (figure on the right). OTT effectively matches digits of the same class while only using supervision from the MNIST dataset. Note that the pixel-level transport plans are both close to the identity meaning that the structure of the images is automatically retrieved. We further illustrate this behaviour by extending this experiment in the supplementary material. From a theoretical perspective, we show that OTT encompasses both Co-OT and GW as well as standard OT. We also show that OTT can be seen as a distance between tensors of any order and thus it can be used to compute tensor barycenters. From an algorithmic point of view, we propose an efficient optimization scheme based on a stochastic mirror descent that allows a drastic reduction of the computational complexity.

The rest of this paper is organized as follows: Section recalls some preliminary knowledge on OT, Co-OT and GW. Section is dedicated to the introduction of our optimal tensor transport (OTT) setting. Section proposes an efficient algorithm for solving OTT. We derive theoretical properties in Section before presenting experimental results on DA and Comparison-based Clustering tasks in Section .

Preliminary Knowledge

In this section, we recall the standard OT (Villani 2008; Peyré, Cuturi et al. 2019), the GW (Memoli 2007; Peyré, Cuturi, and Solomon 2016), and the Co-OT (Redko et al. 2020) formulations. Let p and q be two histograms of respective dimensions I and K . The set of coupling transport plans is defined as $\mathcal{U}_{pq} = \{T \in \mathbb{R}_+^{I \times K} | T \mathbb{1}_K = p, T^\top \mathbb{1}_I = q\}$ where $\mathbb{1}_R$ is a vector of ones of dimension R . The goal in discrete OT is to learn one (in standard OT and GW) or two (in Co-OT) transport plans. Note that for the sake of clarity, we only consider the discrete case here. Nevertheless, all the formulations presented in this section, as well as OTT, can be straightforwardly extended to the continuous case by replacing the sums by integrals over the compared distributions. In this case, the transport plans take the form of joint continuous measures. To prepare for our generalization, we unify the formulations below. In particular, we introduce subscripts and superscripts that are usually not used in the standard formulations. We denote the $(R-1)$ -simplex $\Delta_R = \{(x_r)_{r \in [1, R]} \in \mathbb{R}_+^R | \sum_{r=1}^R x_r = 1\}$.

Optimal Transport (Villani 2008). Let X and Y be two datasets defined over the same feature space \mathcal{X} (e.g. $\mathcal{X} = \mathbb{R}^F$), with respectively $I_1 \in \mathbb{N}$ and $K_1 \in \mathbb{N}$ points with weights $p^1 \in \Delta_{I_1}$ and $q^1 \in \Delta_{K_1}$. The optimal transport plan between X and Y is obtained by solving:

$$\min_{T^1 \in \mathcal{U}_{p^1 q^1}} \sum_{i_1=1}^{I_1} \sum_{k_1=1}^{K_1} \mathcal{L}(X_{i_1}, Y_{k_1}) T_{i_1 k_1}^1 \quad (1)$$

where X_{i_1} is example i_1 in dataset X . Here, \mathcal{L} is a loss function which measures the cost of aligning two examples X_i and Y_k . An extension of OT which is conceptually different from what is covered in this article is the multi-marginal OT (Carlier 2003; Moameni 2014; Pass 2015; Friedland 2020) that aligns $R \geq 3$ datasets simultaneously: \mathcal{L} becomes a function of R parameters and T^1 an R -order tensor.

Co-Optimal Transport (Redko et al. 2020). Co-Optimal Transport also aims at transporting points from two datasets X and Y . However, contrary to standard OT, these datasets may have different feature spaces $\mathcal{X} \subseteq \mathbb{R}^{I_2}$ and $\mathcal{Y} \subseteq \mathbb{R}^{K_2}$ of respective dimensions I_2 and K_2 and equipped with weights $p^2 \in \Delta_{I_2}$ and $q^2 \in \Delta_{K_2}$. The goal is to jointly match the points with a first transport plan T^1 and the features with a second one T^2 . The Co-OT formulation is as follows:

$$\min_{\substack{T^1 \in \mathcal{U}_{p^1 q^1} \\ T^2 \in \mathcal{U}_{p^2 q^2}}} \sum_{i_1, i_2=1}^{I_1, I_2} \sum_{k_1, k_2=1}^{K_1, K_2} \mathcal{L}(X_{i_1 i_2}, Y_{k_1 k_2}) T_{i_1 k_1}^1 T_{i_2 k_2}^2 \quad (2)$$

where $X_{i_1 i_2}$ is the value of feature i_2 for example i_1 .

Gromov Wasserstein (Memoli 2007). Instead of having features describing the examples, let us consider that we only have access to within-dataset pairwise similarities or dissimilarities, that is X and Y are now square matrices of dimensions $I_1 \times I_1$ and $K_1 \times K_1$. It means that the two datasets may have different feature spaces, as in Co-OT, but

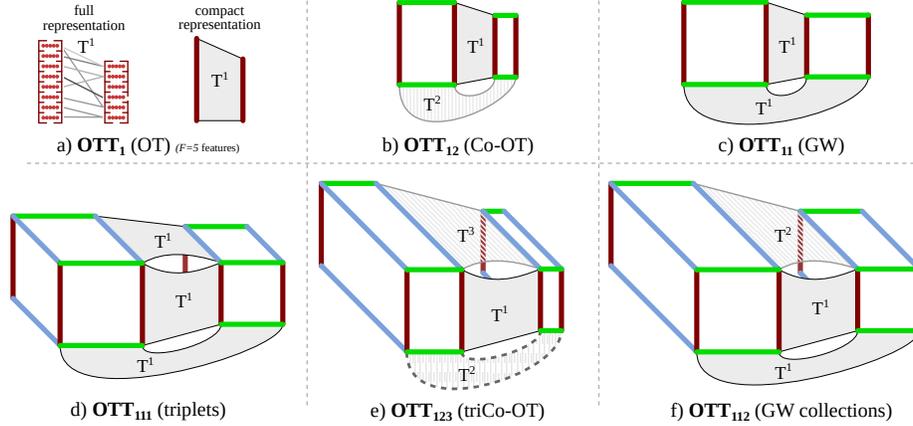


Figure 2: Various formulations of OTT with, each time, the two datasets and the different transport plans (best viewed in color). The subscripts of OTT correspond to the indices of the transport plans used in each dimension.

since these feature spaces are implicit, it is sufficient to learn a single transport plan T^1 . The GW formulation is:

$$\min_{T^1 \in \mathcal{U}_{p^1 q^1}} \sum_{i_1, i_2=1}^{I_1, I_1} \sum_{k_1, k_2=1}^{K_1, K_1} \mathcal{L}(X_{i_1 i_2}, Y_{k_1 k_2}) T^1_{i_1 k_1} T^1_{i_2 k_2} \quad (3)$$

where $X_{i_1 i_2}$ is the (dis)similarity between examples i_1 and i_2 . It is typical, for both Co-OT and GW, to use a comparison/loss function \mathcal{L} (often the squared difference) that operates on two numbers. In Co-OT, \mathcal{L} compares the value of a feature from one point in \mathcal{X} with one feature from a point in \mathcal{Y} . In GW, it compares an entry of the pairwise matrix X to one in Y . Both formulations can be extended by allowing \mathcal{L} to compare more complex entries such as F -dimensional vectors in \mathbb{R}^F . As illustrated in the top row of Figure 2, corresponding to the formulations of Equations (1), (2), and (3), although OT, Co-OT and GW solve different problems, they still share common principles. Below, we propose a new generalized OT formulation that encompasses all of them.

Optimal Tensor Transport (OTT)

Given the notational complexity involved in our generic formulation, let us first explain the intuition behind the subscripts associated with OTT as illustrated in Figure 2. Both Co-OT and GW work on matrices (that is tensors of order $D = 2$) and thus will be represented with 2 digits. Since Co-OT uses $A = 2$ different transport plans T^1 and T^2 , computing Co-OT boils down to solving OTT_{12} as defined below. On the other hand, GW uses the same plan T^1 for both dimensions, thus corresponding to OTT_{11} . Note that dimensions that share a transport plan must have the same sizes. Thus, GW (OTT_{11}) deals with square matrices.

Starting to generalize, when working with tensors of order D , a given OT extension considers $A \leq D$ transport plans and associates a transport plan (index) to each dimension. This is done by specifying an *affectation function* $f : \llbracket 1, D \rrbracket \rightarrow \llbracket 1, A \rrbracket$ or equivalently, a D -tuple of transport plan indices, that is $f \in \llbracket 1, A \rrbracket^D$. For instance, Co-OT uses $f = (1, 2)$ which corresponds to the subscript in OTT_{12} .

For a given $f \in \llbracket 1, A \rrbracket^D$, we can now detail our OTT_f formulation (denoted OTT when no ambiguity arises) that defines a distance between two datasets X and Y , represented as order $D+1$ tensors of respective size $(I_{f(1)} \dots I_{f(D)}, F)$ and $(K_{f(1)} \dots K_{f(D)}, F)$. The first D dimensions will be matched between the two datasets using the transport plans, while the last dimension (F) is the feature dimension used to compare 2 points with the loss \mathcal{L} . To simplify the rest of the paper, we will suppose that $F = 1$, as done in Co-OT and GW above. The OTT distance between X and Y relies on finding a list of optimal transport plans $(T^a)_{a \in \llbracket 1, A \rrbracket}$ under constraints on the marginals defined respectively by the weight vectors $(p^a)_{a \in \llbracket 1, A \rrbracket}$ and $(q^a)_{a \in \llbracket 1, A \rrbracket}$. OTT is defined as:

$$\text{OTT}_f(X, Y, (p^a)_a, (q^a)_a) = \min_{\forall a T^a \in \mathcal{U}_{p^a q^a}} \mathcal{E}_f(X, Y, (T^a)_a) \quad (4)$$

where $\mathcal{E}_f(X, Y, (T^a)_{a \in \llbracket 1, A \rrbracket}) =$

$$\sum_{i_1, \dots, i_D=1}^{I_{f(1)}, \dots, I_{f(D)}} \sum_{k_1, \dots, k_D=1}^{K_{f(1)}, \dots, K_{f(D)}} \mathcal{L}(X_{i_1 \dots i_D}, Y_{k_1 \dots k_D}) \prod_{d=1}^D T^a_{i_d k_d}$$

where $X_{i_1 \dots i_D}$ is the entry at position $i_1 \dots i_D$ in the tensor X .

From this general formulation and looking at Equations 1, 2 and 3 with the support of Figure 2, one can check that OT corresponds to OTT_1 (with F possibly > 1), Co-OT is equivalent to OTT_{12} and GW corresponds to OTT_{11} . Our OTT formulation makes it possible to handle new forms of datasets as illustrated in the second row of Figure 2. In the experiments (see Section), we will specifically consider two versions of OTT, each with order 3 tensors: (i) OTT_{111} corresponds to datasets of triplets (like GW but with triplets instead of pairs); (ii) OTT_{112} works with datasets that are collections of adjacency matrices. Figure 1 gives an illustration of a third kind of datasets, where OTT_{123} has been applied to collections of images, like Co-OT but with three dimensions.

It is worth mentioning that the question that we tackle here is reminiscing of another problem in the literature: the D -regular hypergraphs (Berge 1984) matching. Such a problem is indeed equivalent to OTT in the particular case where all the transport plans are identical. But it uses either a different

Algorithm 1: OTT

Require: datasets X, Y , weights $(p^a)_{a \in \llbracket 1, A \rrbracket}, (q^a)_{a \in \llbracket 1, A \rrbracket}$, loss function \mathcal{L} , nb. of samples M , regularization ϵ

- 1: $\forall a \in \llbracket 1, A \rrbracket, T^a = p^a q^{a\top}$
- 2: **for** $s=0$ to $S-1$ **do**
- 3: **for** $a=1$ to A **do**
- 4: $\widehat{\nabla}_{T^a} \mathcal{E} = M$ gradient samples using Equation (7)
- 5: $T^a = \min_{T \in \mathcal{U}_{p^a q^a}} \langle \widehat{\nabla}_{T^a} \mathcal{E}, T \rangle + \epsilon KL(T, T^a)$
- 6: **end for**
- 7: **end for**

formulation or different constraints on the matching. Zass and Shashua (2008) proposes to find a soft matching between D -regular hypergraphs, with uniform inequality constraints, using a Kullback-Leibler objective function. Duchenne et al. (2011) also matches hypergraphs, with a formulation similar to OTT but uses only row constraints on the matching matrix. Finally, (Peyré et al. 2016) and (Ning and Georgiou 2014) propose to represent examples as PSD matrices and to align those matrices using a single transport plan where each entry is also a PSD matrix instead of a real value.

Algorithm to Solve OTT

In this section, we detail how to efficiently solve the main optimization problem behind Equation 4. The most used method for solving GW is called EGW (Peyré, Cuturi, and Solomon 2016). It can be seen as a Mirror Descent scheme (Beck and Teboulle 2003) with the Kullback-Leibler divergence on a regularized version of GW: $\min_{T \in \mathcal{U}_{p^1 q^1}} \mathcal{E}(T) + KL(T, p^1 q^{1\top})$. The idea of the Mirror Descent algorithm is to interpret the usual gradient descent, at a point x , as a minimization of the sum of a linearization of the desired function $h: \langle \nabla_x h, \bullet \rangle$ plus a regularization term $\epsilon \|x - \bullet\|_2^2$. Instead of using the Euclidean distance, Peyré, Cuturi, and Solomon (2016) use the KL divergence. Thus, at a point T^1 , Peyré, Cuturi, and Solomon (2016) show that the minimization becomes equivalent to the entropy regularized OT problem (Cuturi 2013):

$$\min_{T \in \mathcal{U}_{p^1 q^1}} \langle \nabla_{T^1} \mathcal{E}, T \rangle + \epsilon KL(T, p^1 q^{1\top}). \quad (5)$$

Xu et al. (2019) based on the work of Xie et al. (2020) change the uniform distribution $p^1 q^{1\top}$ in Equation (5) to the previous transport plan T^1 . In fact, this is equivalent to applying a Mirror Descent algorithm on the original GW problem (see Equation (3)) instead of the regularized one (see supplementary material for more details). Thus, to solve the OTT problem, we use a Mirror Descent algorithm with the KL divergence. When the goal is to find multiple transport plans, we propose to use an alternating approach, similar to the Co-OT solver, where each transport plan is optimized in turn while the others remain fixed. In summary, we combine the ideas of existing solvers for Co-OT and GW and apply an alternate Mirror Descent algorithm with the KL divergence for OTT, with the main bottleneck being the computation of the gradient of \mathcal{E} . The pseudo-code of our approach is presented in Algorithm 1. The main steps are the following:

Step 1: We initialize $(T^a)_{a \in \llbracket 1, A \rrbracket}$ with the marginal product.
Step 2: We compute the gradient of \mathcal{E} . For the sake of clarity, we assume that the aligned tensors are ‘‘cubic’’, that is all their dimensions are of the same size N . In this case, the overall gradient with respect to T^a is a N^2 matrix:

$$\nabla_{T^a} \mathcal{E} = \sum_{\{d' | f(d')=a\}} \sum_{\substack{i_1, \dots, i_{d'-1}=1 \\ i_{d'+1}, \dots, i_D=1}} \sum_{\substack{k_1, \dots, k_{d'-1}=1 \\ k_{d'+1}, \dots, k_D=1}} \mathcal{L} \left(\begin{array}{c} X_{i_1 \dots i_{d'-1}, \bullet, i_{d'+1} \dots i_D} \\ Y_{k_1 \dots k_{d'-1}, \bullet, k_{d'+1} \dots k_D} \end{array} \right) \prod_{d=1 | d \neq d'}^D T_{i_d k_d}^{f(d)}. \quad (6)$$

Note that computing the overall gradient exactly would be too expensive. Indeed, a naive approach leads to $O(N^{2D})$ operations which is prohibitively high. To simplify the computation, a first idea would be to generalize the approach used for GW by Peyré, Cuturi, and Solomon (2016) to our problem. This would reduce the complexity to $O(N^{D+1})$ for a particular class of functions \mathcal{L} , notably the square loss. We provide a proof of this approach in the supplementary material. Nevertheless, this remains too expensive as soon as $D = 3$. Thus, instead, we propose to use a stochastic Mirror Descent (Zhou et al. 2017; Zhang and He 2018; Hanzly and Richtárik 2021). This idea was used for the GW problem by Kerdoncuff, Emonet, and Sebban (2021) and we generalize it to our OTT problem. The main idea is to notice that the gradient of \mathcal{E} with respect to T^a can be seen as a sum of expectations over matrices of size N^2 , denoted $(\mathbf{C}^{d'})_{\{d' | f(d')=a\}}$, such that:

$$\mathbb{P} \left(\mathbf{C}^{d'} = \mathcal{L} \left(\begin{array}{c} X_{i_1 \dots i_{d'-1}, \bullet, i_{d'+1} \dots i_D} \\ Y_{k_1 \dots k_{d'-1}, \bullet, k_{d'+1} \dots k_D} \end{array} \right) \right) = \prod_{d=1 | d \neq d'}^D T_{i_d k_d}^{f(d)}$$

with

$$\sum_{\substack{i_1, \dots, i_{d'-1}=1 \\ i_{d'+1}, \dots, i_D=1}} \sum_{\substack{k_1, \dots, k_{d'-1}=1 \\ k_{d'+1}, \dots, k_D=1}} \prod_{\substack{d=1 \\ d \neq d'}}^D T_{i_d k_d}^{f(d)} = 1$$

since $\forall a \in \llbracket 1, A \rrbracket, \sum_{i,k=1}^{I_a, K_a} T_{ik}^a = 1$. The gradient can then be reformulated as:

$$\nabla_{T^a} \mathcal{E} = \sum_{\{d' | f(d')=a\}} \mathbb{E} \left(\mathbf{C}^{d'} \right). \quad (7)$$

It means that one may obtain an unbiased estimate of the gradient in $O(MN^2)$ operations where M is the number of samples to estimate the expectations.

Step 3: The last step (line 5) requires to solve a regularized OT problem, that can be efficiently solved using a Sinkhorn solver (Xu et al. 2019; Cuturi 2013).

We refer the interested reader to Peyré, Cuturi, and Solomon (2016) and Xu et al. (2019) for an analysis of the efficiency of the Mirror Descent algorithm, and to Kerdoncuff, Emonet, and Sebban (2021) for investigations on the precision of the stochastic approximation of the gradient. We also provide in the supplementary material an experiment specific to our new formulation to show how well the gradient is approximated with an increasing order D of the tensors.

Theoretical Results

In this section, we derive two main theoretical results. Theorem 1 shows that as long as the cost function is a proper distance, then OTT is a distance between D -order tensors. Thus, we can naturally define an OTT barycenter between tensors. Theorem 2 states that the optimal barycenter can be found in closed form for particular loss functions.

Theorem 1. *OTT is a distance between weighted tensors $(X, (p^a)_{a \in \llbracket 1, A \rrbracket})$ and $(Y, (q^a)_{a \in \llbracket 1, A \rrbracket})$ represented in canonical form (Definition 1 in the supplementary material), for any affectation function f , as long as \mathcal{L} is a proper distance.*

The proof is provided in the supplementary material. This result notably extends the distance proof of Co-OT (Redko et al. 2020) to matrices of different sizes and to non-uniform weights. Even though their comparison with OTT is out of the scope of this paper, notice that other distances exist between higher-order tensors (De Lathauwer, De Moor, and Vandewalle 2000; Liu, Liu, and Chan 2010; Lai et al. 2013).

Since OTT is a distance, we can define an OTT barycenter between several tensors with any affectation function f .

Definition 1. (*OTT barycenter*) *Assume that we are given $B \geq 1$ weighted tensors of sizes $((K_a^b)_{a \in \llbracket 1, A \rrbracket})_{b \in \llbracket 1, B \rrbracket}$ denoted $(X^b \in \mathbb{R}^{K_{f(1)}^b \dots K_{f(D)}^b}, (q^{a,b} \in \Delta_{K_a^b})_{a \in \llbracket 1, A \rrbracket})_{b \in \llbracket 1, B \rrbracket}$. Let $\lambda \in \Delta_B$ be the weights quantifying the importance of each tensor. For fixed size $(I_a)_{a \in \llbracket 1, A \rrbracket}$ and weights $(p^a \in \Delta_{I_a})_{a \in \llbracket 1, A \rrbracket}$, the OTT barycenter is defined as*

$$\min_{X \in \mathbb{R}^{I_{f(1)} \dots I_{f(D)}}} \sum_{b=1}^B \lambda_b \text{OTT}(X, X^b, (p^a)_a, (q^{a,b})_a). \quad (8)$$

Note that the barycenter could also be defined in a similar manner with the marginals $(p^a)_{a \in \llbracket 1, A \rrbracket}$ not fixed.

To solve Problem (8), we propose to minimize alternatively the objective function w.r.t. X and $(T^{a,b})_{a \in \llbracket 1, A \rrbracket}$, the transport plans between X and X^b . The latter can be found independently for each $b \in \llbracket 1, B \rrbracket$ using Algorithm 1. Interestingly, X can be obtained in closed form for particular loss functions, which generalizes, in particular to Co-OT, a known result for OT and GW (Peyré, Cuturi, and Solomon 2016). This is summarized in the next theorem.

Theorem 2. *Assume that the loss \mathcal{L} is continuous and can be written as $\mathcal{L}(x, y) = f_1(x) + f_2(y) - h_1(x)h_2(y)$ with four functions (f_1, f_2, h_1, h_2) such that $\frac{f_1'}{h_1'}$ is invertible. Further assume that $\mathcal{L}(x, y) \xrightarrow{x \rightarrow \pm\infty} +\infty$. For fixed $((T^{a,b})_{a \in \llbracket 1, A \rrbracket})_{b \in \llbracket 1, B \rrbracket}$, for all $(i_d \in \llbracket 1, I_{f(d)} \rrbracket)_{d \in \llbracket 1, D \rrbracket}$, the optimal solution X_{i_1, \dots, i_D}^* of Problem (8) is equal to*

$$\left(\frac{f_1'}{h_1'} \right)^{-1} \left(\sum_{b=1}^B \lambda_b \sum_{k_1, \dots, k_D=1}^{K_{f(1)}^b, \dots, K_{f(D)}^b} h_2(X_{k_1 \dots k_D}^b) \prod_{d=1}^D \frac{T_{i_d k_d}^{f(d), b}}{p_{i_d}^{f(d)}} \right).$$

In particular, when \mathcal{L} is the squared euclidean distance,

$$X_{i_1, \dots, i_D}^* = \sum_{b=1}^B \lambda_b \sum_{k_1, \dots, k_D=1}^{K_{f(1)}^b, \dots, K_{f(D)}^b} X_{k_1 \dots k_D}^b \prod_{d=1}^D \frac{T_{i_d k_d}^{f(d), b}}{p_{i_d}^{f(d)}}.$$

Note that to obtain a barycenter using loss functions that are not covered by Theorem 2, for example the absolute loss, one can resort to a gradient based optimization scheme.

In the next section, we present experiments focused on 3D-tensors alignments. Nevertheless, it is worth noticing that our theoretical results and Algorithm 1 hold for any tensor order and thus might be used with $D = 4$, for example in comparison based learning tasks (Ghoshdastidar, Perrot, and von Luxburg 2019) or to match hypergraphs (Berge 1984).

Experiments

In this section, we illustrate the interest of OTT on two different tasks¹. First, following the success of OT in Domain Adaptation (Courty et al. 2016), we propose to predict the genres of recent movies based on labeled older movies by relying only on users preferences. We advantageously use a 3D-tensor formulation to take into account the particularity of each user. In a second experiment, we use the OTT barycenter in a Comparison-Based Clustering task.

Domain Adaptation (DA)

We consider a DA task on the Movielens dataset (Harper and Konstan 2015). The goal is to adapt a model learned on old movies (source) to predict the genres of new movies (target).

Datasets. We build two 3-orders tensor X^s (source) and X^t (target) based on the ratings of the users. The entry (i, j, k) in X^s (and similarly for X^t) is 1 if the user i preferred the movie j over the movie k , -1 if the movie k is preferred over the movie j and 0 if the user i cannot choose. As the users did not rate every movie, we use the 0.33 percentile of their personal rates as a default rating. For both the new and old movies, we identify 4 different groups of movies: *Thriller/Crime/Drama (T)*, *Fantasy/Sci-Fi (F)*, *War/Western (W)*, and *Children's/Animation (C)*. We then create 6 pairwise binary classification datasets of 200 movies each by selecting 2 classes among the four aforementioned ones. We assume that we have access to all the labels for the old movies (source) but only to a single random label per class for the new movies (target). The goal is to learn a model that is as accurate as possible on the target. Since many movies have a small number of ratings and many users only rated a few movies, we focus on the 100 users with the highest number of ratings and the 200 most rated films for those users.

Baselines. Even though OTT₁₂₂ is, to the best of our knowledge, the first algorithm that allows direct DA on such tensor-based datasets, we still propose various baselines by reducing the X^s and X^t tensors into matrices by averaging along one dimension. **Rdm** is a first naive baseline that simply outputs random labels. For the next three baselines, we average over the user dimension. Then, **SVM** applies a SVM (Cortes and Vapnik 1995) classifier only on the target domain, using the columns of the matrix as features. **S-GWL** (Xu, Luo, and Carin 2019) interprets the obtained matrices as adjacency matrices of graphs and matches the nodes of the two graphs. **GW** (Peyré, Cuturi, and Solomon

¹The code to reproduce all the experiments is available online: https://github.com/Hv0nnus/Optimal_Tensor_Transport

Datasets	SVM	S-GWL	GW	Co-OT	OTT
T,F	62.5	63.0	62.0	72.0	80.8 ±1
T,C	69.0	77.0	78.0	83.0	97.0 ±0
T,W	32.5	61.0	63.0	65.5	71.3 ±5
F,C	74.5	72.0	74.0	74.0	70.2±4
F,W	53.0	53.0	60.5	47.0	67.9 ±2
C,W	60.0	57.0	52.0	67.5	76.8 ±6
AVG	58.6	63.8	64.9	68.2	77.3 ±3
AVG ^{best}	58.6	66.3	71.0	70.7	78.9 ±3
Time (s)	0.1	94	673	4	5940

Table 1: Accuracy on 6 DA tasks with the hyperparameters found using the unsupervised proposed method. To evaluate the best possible performance reachable by each method, AVG^{best} displays the accuracy with the best hyperparameters using the ground truth of the target domain.

2016) solves the GW problem directly on the obtained matrices. The last baseline, **Co-OT**, uses a matrix obtained by averaging over one movie dimension, which leads to a matrix (users, movies). The two axes are then mapped jointly between the new and old movies. For all the methods that provide a transport plan T between the movies, the class of a target movie y_j^t is predicted via label propagation (Redko et al. 2019a) of the source label y^s , that is $y_k^t = y^s T_{\cdot k}$. The stochastic methods are run 10 times and the mean and standard deviation are reported.

Experimental setup and hyperparameter tuning. As the initialization is key to avoid local minima, we take advantage of both the labels and our stochastic algorithm by sampling only the labelled points in the source and target for the first gradient estimation. The squared euclidean loss is used for \mathcal{L} and we estimate the gradient of OTT using $M = 1000$ samples. S is set to 1000 iterations in Algorithm 1. For each method that uses the OT Sinkhorn solver, notably OTT, we replace it with the semi-supervised algorithm OTDA proposed by Courty et al. (2016) which adds a $l_p - l_1$ regularization to take advantage of the available source labels. In DA, tuning the hyperparameters is often key as there is not enough target labeled movies. As the goal of DA is to reduce the divergence between the two datasets (Ben-David et al. 2007; Redko et al. 2019b), we can use the distance between the source and the target as a criterion to choose the hyperparameters for each method. To compute the OTT distance, we resort to the sampling scheme already used to approximate the GW distance in Kerdoncuff, Emonet, and Sebban (2021). The Kullback-Leibler regularization parameter ϵ of the Sinkhorn method (Cuturi 2013) is selected in the range $[10^{-5}, 10^2]$ and the class regularization η of OTDA (Courty et al. 2016) in $[10^{-4}, 10^1]$. The hyperparameters selection is limited to 24 hours for each method and dataset.

Results. The accuracy of each method is reported in Table 1. OTT achieves better performances than the other baselines on 5 out of 6 datasets. This result was expected as OTT is the only method which takes full advantage of the 3D structure of the data. Interestingly, OTT still behaves better than

the baselines even when one uses the ground truth over the target domain to tune their hyperparameters (that would be cheating) as shown in the line AVG^{best} of Table 1.

We now analyze the impact of the different hyperparameters on the accuracy. We report the results on each dataset in the supplementary material and only consider the global average in Figure 3. The leftmost plot displays the accuracy for increasing values of the KL regularization parameter ϵ . The black markers correspond to the lowest achieved distance for each method. It is worth noting that this usually corresponds to a reasonable accuracy, which supports our hyperparameter tuning procedure. We notice a similar behaviour for the η parameter of OTDA as reported in the supplementary material. In Figure 3 (middle), we report the target accuracy with respect to the number of target labels available. We can notice that OTT is always better, even in the completely unsupervised scenario. Lastly, in the experiments reported in Table 1, we never use the fact that the users comparing the movies are the same for both old and new movies. Here, we study the impact of making this information available. To this end, we fix the transport plan for an increasing number of users. Figure 3 (right) shows that this information greatly improves the target accuracy of the methods that can handle it, especially OTT. Interestingly, as indicated with the black marker, the smallest distance is achieved with the highest number of known pairings, which corresponds to the highest number of constraints on the users transport plan. This supports the key assumption of this experiment: a good matching between users leads to a better matching of similar movies. This also highlights a limit of a mirror descent-based solver as it struggles to find the global minimum without this additional information.

Comparison Based Clustering

In this second series of experiments, we show that OTT barycenters can be used competitively to address an unbalanced comparison-based clustering task.

Comparison-based learning deals with the problem of learning from examples when neither an explicit representation nor a pairwise distance matrix is available (Vikram and Dasgupta 2016; Ukkonen 2017; Emamjomeh-Zadeh and Kempe 2018; Perrot, Esser, and Ghoshdastidar 2020). Instead, it is assumed that only triplet comparisons of the form “example x_i is closer to x_j than to x_k ” are available. This field stems from the fact that relative judgments are usually easier than absolute ones for human observers (Shepard 1962; Young 1987; Stewart, Brown, and Chater 2005). For example, triplet-based queries are easier to answer than exact distance estimations. Given a set of examples and a given number of triplet comparisons, a dataset can be represented as a third order tensor where the entry (i, j, k) contains 1 if example x_i is closer to x_j than to x_k and -1 otherwise. In comparison based clustering, the goal is to identify relevant groups in the examples, using only the information contained in the aforementioned tensor. As the three dimensions of the cubic tensor correspond to the same points we will use the same transport plan for all the dimensions, that is OTT_{111} .

Setting. To show the interest of our method for clustering unbalanced triplet datasets, we take inspiration from the experimental setup of Perrot, Esser, and Ghoshdastidar (2020).

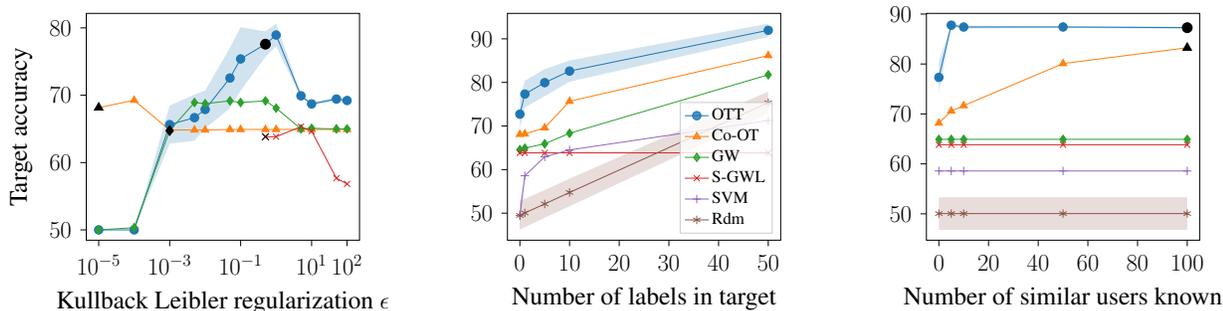


Figure 3: Target accuracy averaged over all the datasets. The shadow area represents the standard deviation for the stochastic methods. (Left) Target accuracy for various values of ϵ . The black symbols correspond to the value of ϵ associated with the lowest distance on average of each method. (Middle) Target accuracy for an increasing target supervision. (Right) Target accuracy for an increasing number of similar known users who rated both old and new movies.

nb. points per class	AddS3	AddS3 _s	t-STE	t-STE _s	OTT
200,20,20	43±12	80±17	56±7	91±4	91±4
30,3,1	28±9	82±17	49±14	91±14	89±14
30,3,3	37±13	78±23	52±09	93±19	87±19
300,30,10	28±4	83±07	48±10	89±4	89±04
AVG	34±9	81±16	51±10	91±8	89±10

Table 2: ARI (in percentage) for unbalanced comparison-based clustering on MNIST. Each line corresponds to the average over 10 combinations of classes, each run 10 times.

For a given dataset, we find the OTT_{111} barycenter ($b = 1$) of size (I_1, I_1, I_1) where I_1 is the number of clusters that we are looking for. The intuition is that similar examples should be sent by the transport plan to the same point in the barycenter since the latter summarizes the initial points.

Datasets. We consider some 3-class unbalanced subsamples of the MNIST dataset (LeCun et al. 1998). For a given number of examples per class (for example, 200,20,20), we consider 10 random draws for the 3 classes and, for each of these, we further consider 10 random draws for the actual images. Given N points in each unbalanced dataset, we randomly select $N \log(N)^3$ triplets of the form $d(x_i, x_j) > d(x_i, x_k)$ as suggested by Perrot, Esser, and Ghoshdastidar (2020). The distance between two digits is the euclidean distance after an UMAP projection in 2 dimensions. To simulate a real dataset, some noise is added by randomly flipping $d(x_i, x_j) > d(x_i, x_k)$ to $d(x_i, x_j) < d(x_i, x_k)$ with probability 0.1 for each triplet selected.

Baselines. We use two main triplet clustering baselines: (i) **t-STE** (Van Der Maaten and Weinberger 2012) which projects the triplets into a vector space followed by k-means (Lloyd 1982), and (ii) **AddS3** (Perrot, Esser, and Ghoshdastidar 2020) which estimates a pairwise similarity matrix also followed by k-means. Moreover, as the OT formulation requires the marginal as a prior, we assume that the proportions of the clusters are known. To stay fair, we propose two variants (AddS3_s, t-STE_s) of the previous baselines

where we replace the k-means step by an OT barycenter step which takes the marginal information into account.

Hyperparameters. We use default hyperparameters, reported in the supplementary material, for t-STE, AddS3, and OTT with the KL regularization parameter set to $\epsilon = 0.1$. To ensure convergence, we also set the number of samples $M = 100$ and the number of iteration $S = 500$ between each of the 20 barycenter updates. Finally, to take advantage of the closed form derived in Theorem 2, we use the squared euclidean loss for OTT.

Results. The Adjusted Rand Index (ARI) (Hubert and Arabie 1985) between the predicted clusters and the ground truth is displayed in Table 2. Overall, OTT has better performances than AddS3_s on average on all datasets while being slightly worse than t-STE_s. Furthermore, for both AddS3 and t-STE, using the unbalancedness information improves the performances. The closeness between our approach and t-STE_s is further investigated in the supplementary material, where we show a theoretical connection between t-STE and the OTT barycenter. The choice of the unbalanced setting is motivated by the fact that the two other baselines do not take into account this information during their first step, while OTT directly uses it in its unique step.

Conclusion

We presented OTT, a new OT formulation that can be used to align high dimensional tensors using potentially several transport plans. OTT generalizes various existing OT problems, such as GW and Co-OT, by defining a new tensor distance. We proposed an efficient algorithm to solve the underlying problem and demonstrated the competitiveness of OTT in DA and Comparison-based clustering. While our new approach unlocks new applications, this comes with a cost. First, despite having access to a solver that drastically reduces the computational complexity of the formulation, it still does not scale well on large datasets with high order tensors. Finally, we leave for future work a natural extension, Fused-OTT, inspired by Vayer et al. (2020), that would combine several OTT problems together. This approach could allow us to align datasets that are independently represented by multiple tensors of potentially different orders.

Acknowledgements

This paper is part of the TADALoT Project funded by the region Auvergne-Rhône-Alpes (France) with the Pack Ambition Recherche (2017, 17 011047 01).

References

- Alvarez-Melis, D.; and Jaakkola, T. 2018. Gromov-Wasserstein Alignment of Word Embedding Spaces. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.
- Arjovsky, M.; Chintala, S.; and Bottou, L. 2017. Wasserstein generative adversarial networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*.
- Beck, A.; and Teboulle, M. 2003. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*.
- Ben-David, S.; Blitzer, J.; Crammer, K.; Pereira, F.; et al. 2007. Analysis of representations for domain adaptation. *Advances in Neural Information Processing Systems*.
- Berge, C. 1984. *Hypergraphs: combinatorics of finite sets*. Elsevier.
- Bunne, C.; Alvarez-Melis, D.; Krause, A.; and Jegelka, S. 2019. Learning Generative Models across Incomparable Spaces. In *International Conference on Machine Learning*.
- Carlier, G. 2003. On a class of multidimensional optimal transportation problems. *Journal of convex analysis*.
- Chowdhury, S.; and Mémoli, F. 2019. The Gromov-Wasserstein distance between networks and stable network invariants. *Information and Inference: A Journal of the IMA*.
- Cortes, C.; and Vapnik, V. 1995. Support-vector networks. *Machine learning*.
- Courty, N.; Flamary, R.; Habrard, A.; and Rakotomamonjy, A. 2017. Joint distribution optimal transportation for domain adaptation. In *Advances in Neural Information Processing Systems*.
- Courty, N.; Flamary, R.; Tuia, D.; and Rakotomamonjy, A. 2016. Optimal transport for domain adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 39(9): 1853–1865.
- Cuturi, M. 2013. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in Neural Information Processing Systems*.
- Damodaran, B. B.; Kellenberger, B.; Flamary, R.; Tuia, D.; and Courty, N. 2018. DeepJDOT: Deep Joint Distribution Optimal Transport for Unsupervised Domain Adaptation. In *European Conference on Computer Vision*. Springer.
- De Lathauwer, L.; De Moor, B.; and Vandewalle, J. 2000. A multilinear singular value decomposition. *SIAM journal on Matrix Analysis and Applications*.
- Duchenne, O.; Bach, F.; Kweon, I.-S.; and Ponce, J. 2011. A tensor-based algorithm for high-order graph matching. *IEEE transactions on pattern analysis and machine intelligence*.
- Emamjomeh-Zadeh, E.; and Kempe, D. 2018. Adaptive Hierarchical Clustering Using Ordinal Queries. In *Symposium on Discrete Algorithms*.
- Friedland, S. 2020. Tensor optimal transport, distance between sets of measures and tensor scaling. *arXiv preprint arXiv:2005.00945*.
- Friedman, J.; Hastie, T.; Tibshirani, R.; et al. 2001. *The elements of statistical learning*. Springer series in statistics New York.
- Ghoshdastidar, D.; Perrot, M.; and von Luxburg, U. 2019. Foundations of Comparison-Based Hierarchical Clustering. In *Advances in Neural Information Processing Systems*.
- Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative adversarial nets. *Advances in Neural Information Processing Systems*.
- Gulrajani, I.; Ahmed, F.; Arjovsky, M.; Dumoulin, V.; and Courville, A. C. 2017. Improved Training of Wasserstein GANs. In *Advances in Neural Information Processing Systems*.
- Hanzely, F.; and Richtárik, P. 2021. Fastest rates for stochastic mirror descent methods. *Computational Optimization and Applications*.
- Harper, F. M.; and Konstan, J. A. 2015. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems*.
- Hubert, L.; and Arabie, P. 1985. Comparing partitions. *Journal of classification*.
- Kerdoncuff, T.; Emonet, R.; and Sebban, M. 2021. Sampled Gromov Wasserstein. *Machine Learning*.
- Lai, Z.; Xu, Y.; Yang, J.; Tang, J.; and Zhang, D. 2013. Sparse tensor discriminant analysis. *IEEE transactions on image processing*.
- LeCun, Y.; Bottou, L.; Bengio, Y.; and Haffner, P. 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*.
- Liu, Y.; Liu, Y.; and Chan, K. C. 2010. Tensor distance based multilinear locality-preserved maximum information embedding. *IEEE Transactions on neural networks*.
- Lloyd, S. 1982. Least squares quantization in PCM. *IEEE transactions on information theory*.
- Memoli, F. 2007. On the use of Gromov-Hausdorff Distances for Shape Comparison. In *Eurographics Symposium on Point-Based Graphics*. The Eurographics Association.
- Mémoli, F. 2011. Gromov-Wasserstein distances and the metric approach to object matching. *Foundations of computational mathematics*.
- Moameni, A. 2014. Multi-marginal Monge-Kantorovich transport problems: A characterization of solutions. *Comptes Rendus Mathématique*.
- Ning, L.; and Georgiou, T. T. 2014. Metrics for matrix-valued measures via test functions. In *53rd IEEE Conference on Decision and Control*, 2642–2647. IEEE.
- Pass, B. 2015. Multi-marginal optimal transport: theory and applications. *ESAIM: Mathematical Modelling and Numerical Analysis*.
- Perrot, M.; Esser, P. M.; and Ghoshdastidar, D. 2020. Near-optimal comparison based clustering. *arXiv preprint arXiv:2010.03918*.

- Peyré, G.; Chizat, L.; Vialard, F.-X.; and Solomon, J. 2016. Quantum optimal transport for tensor field processing. *arXiv preprint arXiv:1612.08731*.
- Peyré, G.; Cuturi, M.; and Solomon, J. 2016. Gromov-wasserstein averaging of kernel and distance matrices. In *International Conference on Machine Learning*.
- Peyré, G.; Cuturi, M.; et al. 2019. Computational optimal transport. *Foundations and Trends® in Machine Learning*.
- Redko, I.; Courty, N.; Flamary, R.; and Tuia, D. 2019a. Optimal transport for multi-source domain adaptation under target shift. In *The 22nd International Conference on Artificial Intelligence and Statistics*. PMLR.
- Redko, I.; Morvant, E.; Habrard, A.; Sebban, M.; and Ben-nani, Y. 2019b. *Advances in domain adaptation theory*. Elsevier.
- Redko, I.; Vayer, T.; Flamary, R.; and Courty, N. 2020. CO-Optimal Transport. In *Advances in Neural Information Processing Systems*.
- Shen, J.; Qu, Y.; Zhang, W.; and Yu, Y. 2018. Wasserstein distance guided representation learning for domain adaptation. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Shepard, R. N. 1962. The analysis of proximities: Multi-dimensional scaling with an unknown distance function. I. *Psychometrika*.
- Stewart, N.; Brown, G. D. A.; and Chater, N. 2005. Absolute identification by relative judgment. *Psychological review*.
- Ukkonen, A. 2017. Crowdsourced correlation clustering with relative distance comparisons. *arXiv preprint arXiv:1709.08459*.
- Van Der Maaten, L.; and Weinberger, K. 2012. Stochastic triplet embedding. In *2012 IEEE International Workshop on Machine Learning for Signal Processing*. IEEE.
- Vayer, T.; Chapel, L.; Flamary, R.; Tavenard, R.; and Courty, N. 2020. Fused Gromov-Wasserstein distance for structured objects. *Algorithms*.
- Vayer, T.; Flamary, R.; Tavenard, R.; Chapel, L.; and Courty, N. 2019. Sliced Gromov-Wasserstein. In *Advances in Neural Information Processing Systems*.
- Vikram, S.; and Dasgupta, S. 2016. Interactive bayesian hierarchical clustering. In *International Conference on Machine Learning*.
- Villani, C. 2008. *Optimal transport: old and new*. Springer Science & Business Media.
- Xie, Y.; Wang, X.; Wang, R.; and Zha, H. 2020. A fast proximal point method for computing exact wasserstein distance. In *Uncertainty in Artificial Intelligence*.
- Xu, H.; Luo, D.; and Carin, L. 2019. Scalable gromov-Wasserstein learning for graph partitioning and matching. In *Advances in Neural Information Processing Systems*.
- Xu, H.; Luo, D.; Zha, H.; and Duke, L. C. 2019. Gromov-Wasserstein Learning for Graph Matching and Node Embedding. In *International Conference on Machine Learning*.
- Yan, Y.; Li, W.; Wu, H.; Min, H.; Tan, M.; and Wu, Q. 2018. Semi-Supervised Optimal Transport for Heterogeneous Domain Adaptation. In *International Joint Conference on Artificial Intelligence*.
- Young, F. W. 1987. *Multidimensional scaling: History, theory, and applications*. Lawrence Erlbaum Associates.
- Zass, R.; and Shashua, A. 2008. Probabilistic graph and hypergraph matching. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE.
- Zhang, S.; and He, N. 2018. On the convergence rate of stochastic mirror descent for nonsmooth nonconvex optimization. *arXiv preprint arXiv:1806.04781*.
- Zhou, Z.; Mertikopoulos, P.; Bambos, N.; Boyd, S.; and Glynn, P. W. 2017. Stochastic mirror descent in variationally coherent optimization problems. *Advances in Neural Information Processing Systems*.