



HAL
open science

Subspace Detours Meet Gromov-Wasserstein

Clément Bonet, Nicolas Courty, François Septier, Lucas Drumetz

► **To cite this version:**

Clément Bonet, Nicolas Courty, François Septier, Lucas Drumetz. Subspace Detours Meet Gromov-Wasserstein. NeurIPS, workshop on Optimal Transport in Machine Learning, Dec 2021, Virtual-only Conference, France. 10.48550/arXiv.2110.10932 . hal-03426813

HAL Id: hal-03426813

<https://hal.science/hal-03426813>

Submitted on 12 Jul 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Subspace Detours Meet Gromov-Wasserstein

Clément Bonet

Univ. Bretagne Sud, LMBA
F-56000 Vannes
clement.bonet@univ-ubs.fr

Nicolas Courty

Univ. Bretagne Sud, IRISA
F-56000 Vannes
nicolas.courty@irisa.fr

François Septier

Univ. Bretagne Sud, LMBA
F-56000 Vannes
francois.septier@univ-ubs.fr

Lucas Drumetz

IMT Atlantique, Lab-STICC
F-29200 Brest
lucas.drumetz@imt-atlantique.fr

Abstract

In the context of optimal transport methods, the *subspace detour* approach was recently presented by Muzellec and Cuturi (2019). It consists in building a nearly optimal transport plan in the measures space from an optimal transport plan in a wisely chosen subspace, onto which the original measures are projected. The contribution of this paper is to extend this category of methods to the Gromov-Wasserstein problem, which is a particular type of transport distance involving the inner geometry of the compared distributions. After deriving the associated formalism and properties, we also discuss a specific cost for which we can show connections with the Knothe-Rosenblatt rearrangement. We finally give an experimental illustration on a shape matching problem.

1 Introduction

Classical optimal transport (OT) has received lots of attention recently, in particular in Machine Learning for tasks such as generative networks (Arjovsky et al., 2017) or domain adaptation (Courty et al., 2016) to name a few. It generally relies on the Wasserstein distance, that builds an optimal coupling between distributions given their geometry. Yet, this metric lacks from potentially important properties, such as translation or rotation invariance, which can be useful when comparing shapes for instance (Mémoli, 2011; Chowdhury et al., 2021), and cannot be used directly whenever the distributions lie in different metric spaces. In order to alleviate those problems, custom solutions have been proposed, such as (Alvarez-Melis et al., 2019; Cai and Lim, 2020).

Apart from these works, another meaningful OT distance to tackle these problems is the Gromov-Wasserstein (GW) distance, originally proposed in Mémoli (2007, 2011). It is a distance between metric spaces and has several appealing properties such as geodesics or invariances (Sturm, 2012). Yet, the price to be paid lies in its computational complexity, which requires to solve a quadratic optimization problem with linear constraints. A recent line of work tends to compute approximations or relaxations of the original problem, in order to spread its use in more data intensive machine learning applications. For example, Peyré et al. use an entropic regularization in order to iterate several Sinkhorn projections (Cuturi, 2013). A related recent method imposes coupling with low-rank constraints (Scetbon et al., 2021). Vayer et al. proposed a sliced approach to approximate Gromov-Wasserstein. Fatras et al. studied an estimator based on mini-batches. In Chowdhury et al. (2021), authors propose to partition the space and to solve the optimal transport problem between a subset of points, before finding a coupling between all the points.

In this work, we study the *subspace detour* approach for Gromov-Wasserstein. This class of method was first proposed for the Wasserstein setting in Muzellec and Cuturi (2019) and consists in choosing the optimal transport plan between projected measures on a subspace, before finding a coupling on the whole space between the original measures using disintegration. Our main contribution is to derive the subspace detours between different subspaces and to apply it for GW costs. We derive some useful properties as well as closed-form solution between Gaussians. Interestingly enough, we also propose a separable quadratic cost for the GW problem that can be related with a triangular coupling, hence bridging the gap with Knothe-Rosenblatt (KR) rearrangements. Illustrations of the method are also given on a shape matching problem.

2 Background

In this section, we introduce all the necessary material to describe the subspace detour approach, from classical optimal transport and its connection to the Knothe-Rosenblatt rearrangement, before defining subspace optimal couplings via the gluing lemma and measure disintegration. Then, we introduce the Gromov-Wasserstein problem for which we will derive the subspace detour in the next sections.

2.1 Classical optimal transport

Let $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$ be two probability measures. The set of couplings between μ and ν is defined as

$$\Pi(\mu, \nu) = \{\gamma \in \mathcal{P}(\mathbb{R}^d \times \mathbb{R}^d), \pi_{\#}^1 \gamma = \mu, \pi_{\#}^2 \gamma = \nu\}$$

where π^1 and π^2 are the projections on the first and second coordinate (*i.e.* $\pi^1(x, y) = x$), and $\#$ is the push forward operator, defined such that

$$\forall A \in \mathcal{B}(\mathbb{R}^d), T_{\#} \mu(A) = \mu(T^{-1}(A)).$$

Kantorovitch problem There exists several types of coupling between probability measures and a non exhaustive list can be found in (Villani, 2008)[Chapter 1]. Among them, the so called optimal coupling is the minimizer of the following Kantorovitch problem:

$$\inf_{\gamma \in \Pi(\mu, \nu)} \int c(x, y) d\gamma(x, y) \quad (1)$$

with c some cost function. When $c(x, y) = \|x - y\|_2^2$, then it defines the Wasserstein distance

$$W_2^2(\mu, \nu) = \inf_{\gamma \in \Pi(\mu, \nu)} \int \|x - y\|_2^2 d\gamma(x, y). \quad (2)$$

The Kantorovitch problem (1) is known to admit a solution when c is nonnegative and lower semi-continuous (Santambrogio, 2015)[Theorem 1.7]. When the optimal coupling is of the form $\gamma = (Id, T)_{\#} \mu$ with T some deterministic map such that $T_{\#} \mu = \nu$, T is called the Monge map.

In one dimension, with μ atomless, the solution to (2) is a deterministic coupling of the form (Santambrogio, 2015)[Theorem 2.5]

$$T = F_{\nu}^{-1} \circ F_{\mu}, \quad (3)$$

where F_{μ} is the cumulative distribution function of μ and F_{ν}^{-1} the quantile function of ν . This map is also known as the increasing rearrangement.

Knothe-Rosenblatt rearrangement Another interesting coupling is the Knothe-Rosenblatt rearrangement, which takes advantage of the increasing rearrangement in one dimension by iterating over the dimension and disintegrating. Concatenating all the increasing rearrangements between the conditional probabilities, we obtain the KR rearrangement, which turns out to be a nondecreasing triangular map (*i.e.* $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$, for all $x \in \mathbb{R}^d$, $T(x) = (T_1(x_1), \dots, T_j(x_1, \dots, x_j), \dots, T_d(x))$) and for all j , T_j is nondecreasing with respect to x_j), and a deterministic coupling (*i.e.* $T_{\#} \mu = \nu$) (Villani, 2008; Santambrogio, 2015; Jaini et al., 2019).

Carlier et al. made a connection between this coupling and optimal transport by showing that it can be obtained as the limit of optimal transport plans for a degenerated cost

$$c_t(x, y) = \sum_{i=1}^d \lambda_i(t)(x_i - y_i)^2,$$

where for all $i \in \{1, \dots, d\}$, $t > 0$, $\lambda_i(t) > 0$ and for all $i \geq 2$, $\frac{\lambda_i(t)}{\lambda_{i-1}(t)} \xrightarrow[t \rightarrow 0]{} 0$. This cost can be recast as in (Bonnotte, 2013) as $c_t(x, y) = (x - y)^T A_t (x - y)$ where $A_t = \text{diag}(\lambda_1(t), \dots, \lambda_d(t))$. This formalizes into the following Theorem:

Theorem 1 (Carlier et al. (2010); Santambrogio (2015)). *Let μ and ν be two absolutely continuous measures on \mathbb{R}^d , with compact supports. Let γ_t be an optimal transport plan for the cost c_t , let T_K be the Knothe-Rosenblatt map between μ and ν , and $\gamma_K = (Id \times T_K)_{\#} \mu$ the associated transport plan. Then, we have $\gamma_t \xrightarrow[t \rightarrow 0]{\mathcal{D}} \gamma_K$. Moreover, if γ_t are induced by transport maps T_t , then T_t converges in $L^2(\mu)$ when t tends to zero to the Knothe-Rosenblatt rearrangement.*

2.2 Subspace detours and disintegration

Muzellec and Cuturi proposed another OT problem by optimizing over the couplings which share a measure on a subspace. More precisely, they defined subspace optimal plans for which the shared measure is the OT plan between projected measures.

Definition 1 (Subspace-Optimal Plans (Muzellec and Cuturi, 2019) Definition 1). *Let $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$ and let $E \subset \mathbb{R}^d$ be a k -dimensional subspace. Let γ_E^* be an OT plan between $\mu_E = \pi_{\#}^E \mu$ and $\nu_E = \pi_{\#}^E \nu$ (with π^E the orthogonal projection on E). Then the set of E -optimal plans between μ and ν is defined as $\Pi_E(\mu, \nu) = \{\gamma \in \Pi(\mu, \nu) \mid (\pi^E, \pi^E)_{\#} \gamma = \gamma_E^*\}$.*

By the Gluing lemma (Villani, 2008), it is possible to construct a coupling $\gamma \in \Pi(\mu, \nu)$ such that $(\pi^E, \pi^E)_{\#} \gamma = \gamma_E^*$. A way to do that is to rely on disintegration.

Disintegration Let (Y, \mathcal{Y}) and (Z, \mathcal{Z}) be measurable spaces, and $(X, \mathcal{X}) = (Y \times Z, \mathcal{Y} \otimes \mathcal{Z})$ the product measurable space. Then, for $\mu \in \mathcal{P}(X)$, we denote $\mu_Y = \pi_{\#}^Y \mu$ and $\mu_Z = \pi_{\#}^Z \mu$ the marginals, where π^Y (respectively π^Z) is the projection on Y (respectively Z). Then, a family $(K(y, \cdot))_{y \in Y}$ is a disintegration of μ if for all $y \in Y$, $K(y, \cdot)$ is a measure on Z , for all $A \in \mathcal{Z}$, $K(\cdot, A)$ is measurable and

$$\forall \phi \in C(X), \int_{Y \times Z} \phi(y, z) d\mu(y, z) = \int_Y \int_Z \phi(y, z) K(y, dz) d\mu_Y(y),$$

where $C(X)$ is the set of continuous functions on X . We can note $\mu = \mu_Y \otimes K$. K is a probability kernel if for all $y \in Y$, $K(y, Z) = 1$. The disintegration of a measure actually corresponds to conditional laws in the context of probabilities. This concept will allow us to obtain measures on the whole space from marginals on subspaces.

In the case where $X = \mathbb{R}^d$, which is the main case of interest in the remainder of the paper, we have existence and uniqueness of the disintegration (see Box 2.2 of Santambrogio (2015) or Chapter 5 of Ambrosio et al. (2008) for the more general case).

Coupling on the whole set Let's note $\mu_{E^\perp|E}$ and $\nu_{E^\perp|E}$ the disintegrated measures on the orthogonal spaces (*i.e.* $\mu = \mu_E \otimes \mu_{E^\perp|E}$ and $\nu = \nu_E \otimes \nu_{E^\perp|E}$). Then, to get a transport plan between the two original measures on the whole space, we can look for another coupling between disintegrated measures $\mu_{E^\perp|E}$ and $\nu_{E^\perp|E}$. In particular, two such couplings are proposed in Muzellec and Cuturi (2019), the Monge-Independent (MI) plan

$$\pi_{\text{MI}} = \gamma_E^* \otimes (\mu_{E^\perp|E} \otimes \nu_{E^\perp|E})$$

where we take the independent coupling between $\mu_{E^\perp|E}(x_E, \cdot)$ and $\nu_{E^\perp|E}(y_E, \cdot)$ for γ_E^* almost every (x_E, y_E) , and the Monge-Knothe (MK) plan

$$\pi_{\text{MK}} = \gamma_E^* \otimes \gamma_{E^\perp|E}^*$$

where $\gamma_{E^\perp|E}^*((x_E, y_E), \cdot)$ is an optimal plan between $\mu_{E^\perp|E}(x_E, \cdot)$ and $\nu_{E^\perp|E}(y_E, \cdot)$ for γ_E^* almost every (x_E, y_E) . Muzellec and Cuturi observed that MI is more adapted to noisy environments since it only computes the OT plan on the subspace. MK is more suited for applications where we want to prioritize some subspace but where all the directions still contain relevant informations.

2.3 Gromov-Wasserstein

Formally, the Gromov-Wasserstein distance allows to compare metric measure spaces (mm-space), triplets (X, d_X, μ_X) and (Y, d_Y, μ_Y) where (X, d_X) and (Y, d_Y) are complete separable metric spaces and μ_X, μ_Y Borel probability measures on X and Y (Sturm, 2012), by computing

$$GW(X, Y) = \inf_{\gamma \in \Pi(\mu_X, \mu_Y)} \iint L(d_X(x, x'), d_Y(y, y')) d\gamma(x, y) d\gamma(x', y')$$

where L is some loss on \mathbb{R} . It has actually been extended to other spaces by replacing the distances by cost functions c_X and c_Y , as *e.g.* in (Chowdhury and Mémoli, 2019). Furthermore, it has many appealing properties such as having invariances (which depend on the costs).

Vayer studied notably this problem in the setting where X and Y are Euclidean spaces, with $L(x, y) = (x - y)^2$ and $c(x, x') = \langle x, x' \rangle$ or $c(x, x') = \|x - x'\|_2^2$. In particular, let $\mu \in \mathcal{P}(\mathbb{R}^p)$, $\nu \in \mathcal{P}(\mathbb{R}^q)$, the inner-GW problem is defined as

$$\text{InnerGW}(\mu, \nu) = \inf_{\gamma \in \Pi(\mu, \nu)} \iint (\langle x, x' \rangle_p - \langle y, y' \rangle_q)^2 d\gamma(x, y) d\gamma(x', y'). \quad (4)$$

For this problem, a closed-form in one dimension can be found:

Theorem 2 (Vayer (2020) Theorem 4.2.4). *Let $\mu, \nu \in \mathcal{P}(\mathbb{R})$, with μ absolutely continuous with respect to the Lebesgue measure. Let $F_\mu^{\nearrow}(x) := F_\mu(x) = \mu(-\infty, x]$ be the cumulative distribution function and $F_\mu^{\searrow}(x) = \mu(-x, +\infty]$ the anti-cumulative distribution function. Let $T_{asc}(x) = F_\nu^{-1}(F_\mu^{\nearrow}(x))$ and $T_{desc}(x) = F_\nu^{-1}(F_\mu^{\searrow}(-x))$. Then, an optimal solution of (4) is achieved either by $\gamma = (Id \times T_{asc})\#\mu$ or by $\gamma = (Id \times T_{desc})\#\mu$.*

3 Subspace detours for GW

In this section, we propose to extend subspace detours from Muzellec and Cuturi (2019) with Gromov-Wasserstein costs. We show that we can even take subspaces of different dimensions, and still obtain a coupling on the whole space using the Independent or the Monge-Knothe coupling. Then, we derive some properties analogously to Muzellec and Cuturi (2019), as well as some closed-form solutions between Gaussians.

3.1 Motivations

First, we adapt the definition of subspace optimal plans for difference subspaces. Indeed, the Gromov-Wasserstein distance being able to compare data on spaces of different dimensions, we can argue that the main information would not be in the same subspace for both datasets. For example, by rotating a dataset, we would change the subspace of interest and most likely lose information as we can see on Figure 1. On this illustration, we use as a source one moon of the Two moons dataset, and obtain a target by rotating it by an angle of $\frac{\pi}{2}$. As GW with $c(x, x') = \|x - x'\|_2^2$ is invariant with respect to isometries, we are able to recover the exact correspondence between the points. However, when choosing a subspace to project both the source and target, we completely lose the optimal coupling between them. Nonetheless, by choosing more wisely one subspace by dataset (using here the first component of the principal component analysis (PCA) decomposition), we find the right coupling. This illustration underlines the idea that the choice of both subspaces is important. A way of choosing the subspaces could be to project on the subspace containing the more information for each dataset using *e.g.* PCA independently on each distribution. Muzellec and Cuturi proposed to optimize the optimal transport cost with respect to an orthonormal matrix with a projected gradient descent, which could be extended to an optimization over two orthonormal matrices in our context.

By allowing to have different subspaces, we get the following definition of subspace optimal plans.

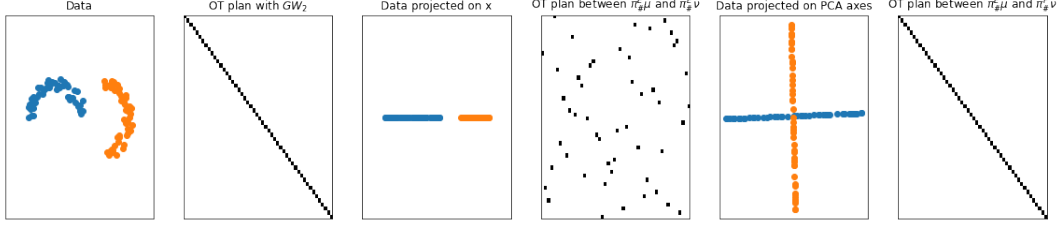


Figure 1: From left to right: Data (moons), OT plan obtained with GW for $c(x, x') = \|x - x'\|_2^2$, Data projected on the 1st axis, OT plan obtained between the projected measures, Data projected on their 1st PCA component, OT plan obtained between the the projected measures

Definition 2. Let $\mu \in \mathcal{P}_2(\mathbb{R}^p)$, $\nu \in \mathcal{P}_2(\mathbb{R}^q)$, E be a k -dimensional subspace of \mathbb{R}^p and F a k' -dimensional subspace of \mathbb{R}^q . Let $\gamma_{E,F}^*$ be an optimal transport plan for GW between $\mu_E = \pi_{\#}^E \mu$ and $\nu_F = \pi_{\#}^F \nu$ (with π^E (resp. π^F) the orthogonal projection on E (resp. F)). Then the set of (E, F) -optimal plans between μ and ν is defined as $\Pi_{E,F}(\mu, \nu) = \{\gamma \in \Pi(\mu, \nu) \mid (\pi^E, \pi^F)_{\#} \gamma = \gamma_{E,F}^*\}$.

Analogously to Muzellec and Cuturi (2019) (Section 2.2), we can obtain from $\gamma_{E,F}^*$ a coupling on the whole set by either defining the Monge-Independent plan $\pi_{\text{MI}} = \gamma_{E,F}^* \otimes (\mu_{E^\perp|E} \otimes \nu_{F^\perp|F})$ or the Monge-Knothe plan $\pi_{\text{MK}} = \gamma_{E,F}^* \otimes \gamma_{E^\perp \times F^\perp|E \times F}^*$ where OT plans are taken with some OT cost such as e.g. GW.

3.2 Properties

Following Muzellec and Cuturi (2019), the Monge-Knothe coupling is the optimal measure among the subspace optimal plans for the corresponding cost. We show it for the Gromov-Wasserstein distance with cost L , which is a direct transposition of Proposition 1 in Muzellec and Cuturi (2019).

Proposition 1. Let $\mu \in \mathcal{P}(\mathbb{R}^p)$ and $\nu \in \mathcal{P}(\mathbb{R}^q)$, $E \subset \mathbb{R}^p$, $F \subset \mathbb{R}^q$, $\pi_{\text{MK}} = \gamma_{E,F}^* \otimes \gamma_{E^\perp \times F^\perp|E \times F}^*$ where γ^* are optimal for the Gromov-Wasserstein problem with cost L . Then we have:

$$\pi_{\text{MK}} \in \underset{\gamma \in \Pi_{E,F}(\mu, \nu)}{\text{argmin}} \iint L(x, x', y, y') d\gamma(x, y) d\gamma(x', y').$$

Key properties of GW that we would like to keep are its invariances. We show in two particular cases that we conserve them on the orthogonal spaces (since the measure on $E \times F$ is fixed).

Proposition 2. Let $\mu \in \mathcal{P}(\mathbb{R}^p)$, $\nu \in \mathcal{P}(\mathbb{R}^q)$, and denote

$$GW_{E,F}(\mu, \nu) = \inf_{\gamma \in \Pi_{E,F}(\mu, \nu)} \iint L(x, x', y, y') d\gamma(x, y) d\gamma(x', y').$$

For $L(x, x', y, y') = (\|x - x'\|_2^2 - \|y - y'\|_2^2)^2$, $GW_{E,F}$ is invariant with respect to translations and isometries on E^\perp and F^\perp .

For $L(x, x', y, y') = (\langle x, x' \rangle_p - \langle y, y' \rangle_q)^2$, $GW_{E,F}$ is invariant with respect to isometries on E^\perp and F^\perp .

We refer to Appendix A.1 for the proofs of the two previous propositions.

3.3 Closed-form between Gaussians

We can also derive explicit formulas between Gaussians in particular cases. Let $q \leq p$, $\mu = \mathcal{N}(m_\mu, \Sigma) \in \mathcal{P}(\mathbb{R}^p)$, $\nu = \mathcal{N}(m_\nu, \Lambda) \in \mathcal{P}(\mathbb{R}^q)$ two Gaussian measures with $\Sigma = P_\mu D_\mu P_\mu^T$ and $\Lambda = P_\nu D_\nu P_\nu^T$. As previously, let $E \subset \mathbb{R}^p$ and $F \subset \mathbb{R}^q$ be respectively k and k' dimensional subspaces. Following Muzellec and Cuturi (2019), we represent Σ in an orthonormal basis of $E \oplus E^\perp$, and Λ in an orthonormal basis of $F \oplus F^\perp$, i.e. $\Sigma = \begin{pmatrix} \Sigma_E & \Sigma_{EE^\perp} \\ \Sigma_{E^\perp E} & \Sigma_{E^\perp} \end{pmatrix}$. Now, let's denote

$$\Sigma/\Sigma_E = \Sigma_{E^\perp} - \Sigma_{EE^\perp}^T \Sigma_E^{-1} \Sigma_{EE^\perp}$$

the Schur complement of Σ with respect to Σ_E . We know that the conditionals of Gaussians are Gaussians, and of covariance the Schur complement (see *e.g.* Rasmussen (2003); Von Mises (1964)).

For $L(x, x', y, y') = (\|x - x'\|_2^2 - \|y - y'\|_2^2)^2$, we have for now no certainty that the optimal transport plan is Gaussian. By restricting the minimization problem to Gaussian couplings, Salmona et al. showed that there is a solution $\gamma^* = (Id, T)_{\#}\mu \in \Pi(\mu, \nu)$ with $\mu = \mathcal{N}(m_\mu, \Sigma)$, $\nu = \mathcal{N}(m_\nu, \Lambda)$ and

$$\forall x \in \mathbb{R}^d, T(x) = m_\nu + P_\nu A P_\mu^T (x - m_\mu) \quad (5)$$

where $A = \left(\tilde{I}_q D_\nu^{\frac{1}{2}} (D_\mu^{(q)})^{-\frac{1}{2}} \quad 0_{q, p-q} \right) \in \mathbb{R}^{q \times p}$ and \tilde{I}_q is of the form $\text{diag}((\pm 1)_{i \leq q})$.

By combining the results of Muzellec and Cuturi (2019) and Salmona et al. (2021), we get the following closed-form for Monge-Knothe couplings.

Proposition 3. *Suppose $p \geq q$ and $k = k'$. For the Gaussian restricted GW problem, a Monge-Knothe transport map between $\mu = \mathcal{N}(m_\mu, \Sigma) \in \mathcal{P}(\mathbb{R}^p)$ and $\nu = \mathcal{N}(m_\nu, \Lambda) \in \mathcal{P}(\mathbb{R}^q)$ is, for all $x \in \mathbb{R}^p$, $T_{\text{MK}}(x) = m_\nu + B(x - m_\mu)$ where*

$$B = \begin{pmatrix} T_{E,F} & 0 \\ C & T_{E^\perp, F^\perp | E, F} \end{pmatrix}$$

with $T_{E,F}$ an optimal transport map between $\mathcal{N}(0_E, \Sigma_E)$ and $\mathcal{N}(0_F, \Lambda_F)$ (of the form (5)), $T_{E^\perp, F^\perp | E, F}$ an optimal transport map between $\mathcal{N}(0_{E^\perp}, \Sigma/\Sigma_E)$ and $\mathcal{N}(0_{F^\perp}, \Lambda/\Lambda_F)$ and C satisfying

$$C = (\Lambda_{F^\perp F} (T_{E,F}^T)^{-1} - T_{E^\perp, F^\perp | E, F} \Sigma_{E^\perp E}) \Sigma_E^{-1}.$$

Proof. See Appendix A.2.1. □

Suppose that $k \geq k'$, $m_\mu = 0$, $m_\nu = 0$ and let $T_{E,F}$ be an optimal transport map between μ_E and ν_F (of the form (5)). We can derive a formula for the Monge-Independent coupling for the inner-GW problem and the Gaussian restricted GW problem.

Proposition 4. $\pi_{\text{MI}} = \mathcal{N}(0_{p+q}, \Gamma)$ where $\Gamma = \begin{pmatrix} \Sigma & C \\ C^T & \Lambda \end{pmatrix}$ with

$$C = (V_E \Sigma_E + V_{E^\perp} \Sigma_{E^\perp E}) T_{E,F}^T (V_F^T + \Lambda_F^{-1} \Lambda_{F^\perp F}^T V_{F^\perp}^T)$$

where $T_{E,F}$ is an optimal transport map, either for the inner-GW problem or the Gaussian restricted problem.

Proof. See Appendix A.2.2. □

3.4 Limit of optimal transport plans?

Another interesting property derived in Muzellec and Cuturi (2019) of the Monge-Knothe coupling is that it can be obtained as the limit of classic optimal transport plans, similar to Theorem 1, using a separable cost of the form

$$c_t(x, y) = (x - y)^T P_t (x - y)$$

with $P_t = V_E V_E^T + t V_{E^\perp} V_{E^\perp}^T$ and (V_E, V_{E^\perp}) an orthonormal basis of \mathbb{R}^p .

However, this property is not valid for the classical Gromov-Wasserstein cost (*e.g.* $L(x, x', y, y') = (d_X(x, x')^2 - d_Y(y, y')^2)^2$ or $L(x, x', y, y') = (\langle x, x' \rangle_p - \langle y, y' \rangle_q)^2$) as the cost is not separable. Motivated by this question, we ask ourselves in the following if we can derive a quadratic optimal transport cost for which we would have this property.

Construction and properties of the Hadamard-Wasserstein problem The main idea of the proof of Theorem 1 in Carlier et al. (2010) is to decompose the objective function as

$$\int c_t(x, y) d\gamma(x, y) = \lambda_1(t) \left(\int (x_1 - y_1)^2 d\gamma(x, y) + \int \sum_{k=2}^d \frac{\lambda_k(t)}{\lambda_1(t)} (x_k - y_k)^2 d\gamma(x, y) \right),$$

before taking the limit $t \rightarrow 0$ which makes the right-hand term vanish and allows to conclude on the limit of the first marginal of the optimal map. Reasoning by induction on the dimension, Carlier et al. are able to deal with one term at a time, and finally show that the limit of the optimal map is the Knothe-Rosenblatt transport (2.1). Another key ingredient is to have access to a unique transport map between measures in \mathbb{R} , as it is the case for the Wasserstein distance with cost $c(x, y) = \frac{1}{2}(x - y)^2$, the Monge map being the increasing rearrangement (3) (it can actually be extended to smoothly strictly convex costs, see Santambrogio (2015)[Theorem 2.9]).

For now, the only cost for which we have an optimal transport map in 1D is for the inner product (Vayer, 2020). Hence, we need a cost which reduces to inner-GW (4) in 1D. A natural choice is therefore to use the following cost:

$$\forall x, x', y, y' \in \mathbb{R}^d, L(x, x', y, y') = \sum_{k=1}^d (x_k x'_k - y_k y'_k)^2 = \|x \odot x' - y \odot y'\|_2^2 \quad (6)$$

as a loss function, where \odot is the Hadamard product (element wise product). We define the following ‘‘Hadamard Wasserstein’’ problem

$$\mathcal{HW}(\mu, \nu) = \inf_{\gamma \in \Pi(\mu, \nu)} \iint \|x \odot x' - y \odot y'\|_2^2 d\gamma(x, y) d\gamma(x', y'). \quad (7)$$

Properties The loss L (6) satisfies well the separability condition and reduces to the inner-GW loss in 1D. We can therefore define a degenerated version of it,

$$\begin{aligned} \forall x, x', y, y', L_t(x, x', y, y') &= \sum_{k=1}^d \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 \\ &= (x \odot x' - y \odot y') A_t (x \odot x' - y \odot y') \end{aligned} \quad (8)$$

with $A_t = \text{diag}(1, \lambda_t^{(1)}, \lambda_t^{(1)} \lambda_t^{(2)}, \dots, \prod_{i=1}^{d-1} \lambda_t^{(i)})$, and such as for all $t > 0$, and for all $i \in \{1, \dots, d-1\}$, $\lambda_t^{(i)} > 0$ and $\lambda_t^{(i)} \xrightarrow[t \rightarrow 0]{} 0$. We denote \mathcal{HW}_t the problem (7) with the degenerate cost (8). We will derive some useful properties which are usual for the regular Gromov-Wasserstein cost.

Proposition 5. *Let $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$.*

1. *The problem (7) always admits a minimizer.*
2. *\mathcal{HW} is a pseudometric (i.e. it is symmetric, nonnegative, $\mathcal{HW}(\mu, \mu) = 0$ and it satisfies the triangle inequality).*
3. *\mathcal{HW} is invariant to reflexion with respect to axes.*

Proof. See Appendix B.1. □

\mathcal{HW} loses some properties compared to GW . Indeed, it is only invariant with respect to axes and it can compare only measures lying in the same Euclidean space in order for the distance to be well defined. Nonetheless, we show in the following that we can derive some links with triangular couplings in the same way as the Wasserstein distance and KR.

We first define a triangular coupling different from the Knothe-Rosenblatt rearrangement in the sense that each map will not be nondecreasing. Indeed, following Theorem 2, the solution of each 1D problem

$$\text{argmin}_{\gamma \in \Pi(\mu, \nu)} \iint (xx' - yy')^2 d\gamma(x, y) d\gamma(x', y')$$

is either $(Id \times T_{\text{asc}})_{\#} \mu$ or $(Id \times T_{\text{desc}})_{\#} \mu$. Hence, at each step $k \geq 1$, if we disintegrate the joint law of the k first variables as $\mu^{1:k} = \mu^{1:k-1} \otimes \mu^{k|1:k-1}$, the optimal transport map $T(\cdot | x_1, \dots, x_k)$ will be the solution of

$$\text{argmin}_{T \in \{T_{\text{asc}}, T_{\text{desc}}\}} \iint (x_k x'_k - T(x_k) T(x'_k))^2 \mu^{k|1:k-1}(dx_k | x_{1:k-1}) \mu^{k|1:k-1}(dx'_k | x'_{1:k-1}).$$

We now state the main theorem where we show that the limit of the OT plans obtained with the degenerated cost will be the triangular coupling we just defined.

Theorem 3. *Let μ and ν be two absolutely continuous measures on \mathbb{R}^d such that $\int \|x\|_2^4 \mu(dx) < +\infty$, $\int \|y\|_2^4 \nu(dy) < +\infty$ and with compact support. Let γ_t be an optimal transport plan for \mathcal{HW}_t , let T_K be the alternate Knothe-Rosenblatt map between μ and ν as defined in the last paragraph, and let $\gamma_K = (Id \times T_K)_\# \mu$ be the associated transport plan. Then, we have $\gamma_t \xrightarrow[t \rightarrow 0]{\mathcal{D}} \gamma_K$. Moreover, if γ_t are induced by transport maps T_t , then $T_t \xrightarrow[t \rightarrow 0]{L^2(\mu)} T_K$.*

Proof. See appendix B.2. □

We report in Appendix C how to compute \mathcal{HW} (7) in the discrete setting.

4 Illustrations

We use the Python Optimal Transport (POT) library (Flamary et al., 2021) to compute the different optimal transport problems involved in this illustration. We are interested here in solving a 3D mesh registration problem, which is a natural application of Gromov-Wasserstein (Mémoli, 2011) since it enjoys invariances with respect to isometries such as permutations, and can also naturally exploit the topology of the meshes. For this purpose, we selected two base meshes from the FAUST dataset (Bogo et al., 2014), which provides ground truth correspondences between shapes. The information available from those meshes are geometrical (6890 vertices positions) and topological (mesh connectivity). These two meshes are represented, along with the visual results of the registration, in Figure 2. In order to visually depict the quality of the assignment induced by the transport map, we propagate through it a color code of the source vertices toward their associated counterpart vertices in the target mesh. Both original color coded source and associated target ground truth are available on the first line of the illustration. To compute our method, we simply use as a natural subspace for both meshes the algebraic connectivity of the mesh topological information, also known as the Fiedler vector (eigenvector associated to the second smallest eigenvalue of the un-normalized Laplacian matrix). Reduced to a 1D optimal transport problem, following Eq. 4, the computation time is very low (~ 5 secs. on a standard laptop), and the associated matching is very good with more than 98% of correct assignments. We qualitatively compare this result to Gromov-Wasserstein mappings induced by different cost functions, in the second line of Figure 2: adjacency (Xu et al., 2019), weighted adjacency (weights are given by distances between vertices), heat kernel (derived from the un-normalized Laplacian) (Chowdhury and Needham, 2021) and finally geodesic distances over the meshes. In average, computing the Gromov-Wasserstein mapping using POT took around 10 minutes of time. Both methods based on adjacency fail to recover a meaningful mapping. Heat kernel allows to map continuous areas of the source mesh, but fails in recovering a global structure. Finally, the geodesic distance gives a much more coherent mapping, but has inverted left and right of the human figure. Notably, a significant extra computation time was induced by the computation of the geodesic distances (~ 1 h/mesh using the NetworkX (Hagberg et al., 2008) shortest path procedure). As a conclusion, and despite the simplification of the original problem, our method performs best, with a speed-up of two-orders of magnitude.

5 Discussion

We proposed in this work to extend the subspace detour approach to different subspaces, and to other optimal transport costs such as Gromov-Wasserstein. Being able to project on different subspaces can be useful when the data are not aligned and do not share the same axes of interest, as well as when we are working between different metric spaces as it is the case for example with graphs. However, a question arising is how to choose these subspaces. Since the method is mostly interesting when we choose one dimensional subspaces, we proposed to use a PCA and to project on the first directions for data embedded in euclidean spaces. For more complicated data such as graphs, we projected onto the Fiedler vector and obtained good results in an efficient way on a 3D mesh registration problem. More generally, Muzellec and Cuturi proposed to perform a gradient descent on the loss with respect to orthonormal matrices. This approach is non-convex and only guaranteed to converge to a local minimum. Designing such an algorithm, which would minimize alternatively between two transformations in the Stiefel manifold, is left for future works.

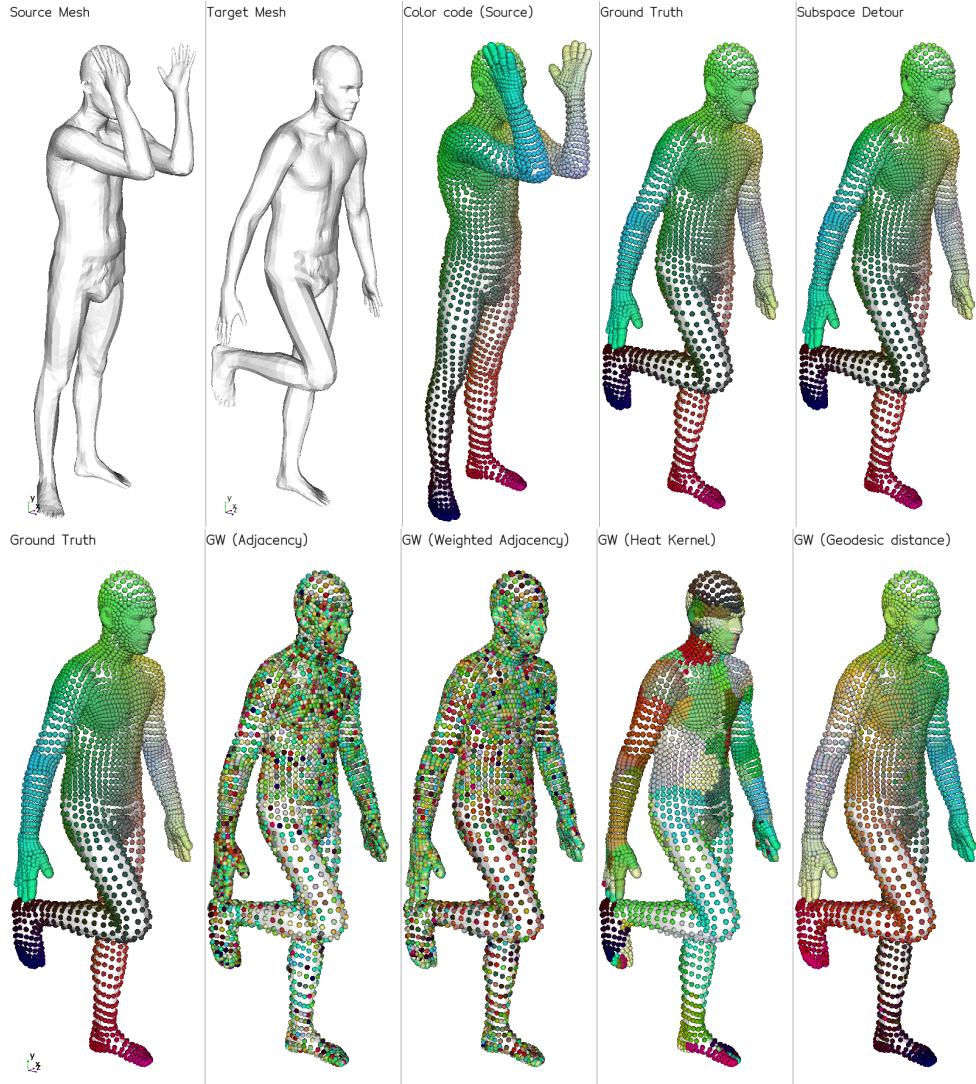


Figure 2: 3D Mesh registration. (First row) source and target meshes, color code of the source, ground truth color code on the target, result of subspace detour using Fiedler vectors as subspace. (Second row) After recalling the expected ground truth for ease of comparison, we present results of different Gromov-Wasserstein mappings obtained with metrics based on adjacency, heat kernel and geodesic distances.

The subspace detour approach for transport problem is meaningful whenever one can identify subspaces that gather most of the information from the original distributions, while making the estimate more robust and with a better sample complexity as far as dimensions are lower. On the computational complexity side, and when we have only access to discrete data, the subspace detour approach brings better computational complexity solely when the subspaces are chosen as one dimensional. Indeed, otherwise, we have the same complexity for solving the subspace detour and solving directly the OT problem (since the complexity only depends on the number of samples). In this case, the 1D projection often gives distinct values for all the samples (for continuous valued data) and hence the Monge-Knothe coupling is exactly the coupling in 1D. As such, information is lost on the orthogonal spaces. It can be artificially recovered by quantizing the 1D values (as experimented in practice in Muzellec and Cuturi (2019)), but the added value is not clear and deserves broader studies. If given absolutely continuous distributions *wrt.* the Lebesgue measure however, this limit does not exist, but comes with the extra cost of being able to compute efficiently the projected measure onto

the subspace, which might require discretization of the space and is therefore not practical in high dimensions.

We also proposed a new quadratic cost \mathcal{HW} that we call Hadamard-Wasserstein, which allows to define a degenerated cost for which the optimal transport plan converges to a triangular coupling. However, this cost loses many properties compared to W_2 or GW , for which we are inclined to use these problems. Indeed, while \mathcal{HW} is a quadratic cost, it uses an euclidean norm between the Hadamard product of vectors and requires the two spaces to be the same (in order to have the distance well defined). A work around in the case $X = \mathbb{R}^p$ and $Y = \mathbb{R}^q$ with $p \leq q$ would be to “lift” the vectors in \mathbb{R}^p into vectors in \mathbb{R}^q with padding as it is proposed in Vayer et al. (2019b), or to project the vectors in \mathbb{R}^q on \mathbb{R}^p as in Cai and Lim (2020). Yet for some applications where only the distance/similarity matrices are available, a different strategy still needs to be found. Another concern is the limited invariance properties (only with respect to axial symmetry symmetry in our case). Nevertheless, we expect that such a cost can be of interest in cases where invariance to symmetry is a desired property, such as in (Nagar and Raman, 2019).

References

- David Alvarez-Melis, Stefanie Jegelka, and Tommi S Jaakkola. Towards optimal transport with global invariances. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1870–1879. PMLR, 2019.
- Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media, 2008.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- Patrick Billingsley. *Convergence of probability measures*. John Wiley & Sons, 2013.
- Federica Bogo, Javier Romero, Matthew Loper, and Michael J Black. Faust: Dataset and evaluation for 3d mesh registration. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3794–3801, 2014.
- Nicolas Bonnotte. *Unidimensional and evolution methods for optimal transportation*. PhD thesis, Paris 11, 2013.
- Yuhang Cai and Lek-Heng Lim. Distances between probability distributions of different dimensions. *arXiv preprint arXiv:2011.00629*, 2020.
- Guillaume Carlier, Alfred Galichon, and Filippo Santambrogio. From knothe’s transport to brenier’s map and a continuation method for optimal transport. *SIAM Journal on Mathematical Analysis*, 41(6):2554–2576, 2010.
- Samir Chowdhury and Facundo Mémoli. The gromov–wasserstein distance between networks and stable network invariants. *Information and Inference: A Journal of the IMA*, 8(4):757–787, 2019.
- Samir Chowdhury and Tom Needham. Generalized spectral clustering via gromov-wasserstein learning. In *International Conference on Artificial Intelligence and Statistics*, pages 712–720. PMLR, 2021.
- Samir Chowdhury, David Miller, and Tom Needham. Quantized gromov-wasserstein, 2021.
- 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. *Advances in neural information processing systems*, 26:2292–2300, 2013.
- Kilian Fatras, Younes Zine, Szymon Majewski, Rémi Flamary, Rémi Gribonval, and Nicolas Courty. Minibatch optimal transport distances; analysis and applications. *arXiv preprint arXiv:2101.01792*, 2021.

- Rémi Flamary, Nicolas Courty, Alexandre Gramfort, Mokhtar Z. Alaya, Aurélie Boisbunon, Stanislas Chambon, Laetitia Chapel, Adrien Corenflos, Kilian Fatras, Nemo Fournier, Léo Gautheron, Nathalie T.H. Gayraud, Hicham Janati, Alain Rakotomamonjy, Ievgen Redko, Antoine Rolet, Antony Schutz, Vivien Seguy, Danica J. Sutherland, Romain Tavenard, Alexander Tong, and Titouan Vayer. Pot: Python optimal transport. *Journal of Machine Learning Research*, 22(78):1–8, 2021. URL <http://jmlr.org/papers/v22/20-451.html>.
- Aric Hagberg, Pieter Swart, and Daniel S Chult. Exploring network structure, dynamics, and function using networkx. Technical report, Los Alamos National Lab.(LANL), Los Alamos, NM (United States), 2008.
- Priyank Jaini, Kira A Selby, and Yaoliang Yu. Sum-of-squares polynomial flow. In *International Conference on Machine Learning*, pages 3009–3018. PMLR, 2019.
- Simon Lacoste-Julien. Convergence rate of frank-wolfe for non-convex objectives. *arXiv preprint arXiv:1607.00345*, 2016.
- Facundo Mémoli. On the use of gromov-hausdorff distances for shape comparison. 2007.
- Facundo Mémoli. Gromov–wasserstein distances and the metric approach to object matching. *Foundations of computational mathematics*, 11(4):417–487, 2011.
- Boris Muzellec and Marco Cuturi. Subspace detours: Building transport plans that are optimal on subspace projections. *arXiv preprint arXiv:1905.10099*, 2019.
- Rajendra Nagar and Shanmuganathan Raman. Detecting approximate reflection symmetry in a point set using optimization on manifold. *IEEE Transactions on Signal Processing*, 67(6):1582–1595, 2019.
- François-Pierre Paty and Marco Cuturi. Subspace robust wasserstein distances. In *International Conference on Machine Learning*, pages 5072–5081. PMLR, 2019.
- Gabriel Peyré, Marco Cuturi, and Justin Solomon. Gromov-wasserstein averaging of kernel and distance matrices. In *International Conference on Machine Learning*, pages 2664–2672. PMLR, 2016.
- Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport: With applications to data science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- Carl Edward Rasmussen. Gaussian processes in machine learning. In *Summer school on machine learning*, pages 63–71. Springer, 2003.
- Antoine Salmona, Julie Delon, and Agnès Desolneux. Gromov-wasserstein distances between gaussian distributions. 2021.
- Filippo Santambrogio. Optimal transport for applied mathematicians. *Birkäuser, NY*, 55(58-63):94, 2015.
- Meyer Scetbon, Gabriel Peyré, and Marco Cuturi. Linear-time gromov wasserstein distances using low rank couplings and costs. *arXiv preprint arXiv:2106.01128*, 2021.
- Karl-Theodor Sturm. The space of spaces: curvature bounds and gradient flows on the space of metric measure spaces. *arXiv preprint arXiv:1208.0434*, 2012.
- Titouan Vayer. *A contribution to Optimal Transport on incomparable spaces*. PhD thesis, 2020.
- Titouan Vayer, Nicolas Courty, Romain Tavenard, and Rémi Flamary. Optimal transport for structured data with application on graphs. In *International Conference on Machine Learning*, pages 6275–6284. PMLR, 2019a.
- Titouan Vayer, Rémi Flamary, Nicolas Courty, Romain Tavenard, and Laetitia Chapel. Sliced gromov-wasserstein. *Advances in Neural Information Processing Systems*, 32:14753–14763, 2019b.

Cédric Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.

Richard Von Mises. *Mathematical theory of probability and statistics*. Academic Press, 1964.

Hongteng Xu, Dixin Luo, and Lawrence Carin. Scalable gromov-wasserstein learning for graph partitioning and matching. *Advances in neural information processing systems*, 32:3052–3062, 2019.

A Subspace detours

A.1 Proofs

Proof of Proposition 1. $\forall \gamma \in \Pi_{E,F}(\mu, \nu)$,

$$\begin{aligned} & \iint L(x, x', y, y') d\gamma(x, y) d\gamma(x', y') \\ &= \iint \left(\iint L(x, x', y, y') \gamma_{E^\perp \times F^\perp | E \times F}((x_E, y_F), (dx_{E^\perp}, dy_{F^\perp})) \gamma_{E^\perp \times F^\perp | E \times F}((x'_E, y'_F), (dx'_{E^\perp}, dy'_{F^\perp})) \right) \\ & d\gamma_{E \times F}^*(x_E, y_F) d\gamma_{E \times F}^*(x'_E, y'_F) \end{aligned}$$

However, for $\gamma_{E \times F}^*$ a.e. $(x_E, y_F), (x'_E, y'_F)$,

$$\begin{aligned} & \iint L(x, x', y, y') \gamma_{E^\perp \times F^\perp | E \times F}((x_E, y_F), (dx_{E^\perp}, dy_{F^\perp})) \gamma_{E^\perp \times F^\perp | E \times F}((x'_E, y'_F), (dx'_{E^\perp}, dy'_{F^\perp})) \\ & \geq \iint L(x, x', y, y') \gamma_{E^\perp \times F^\perp | E \times F}^*((x_E, y_F), (dx_{E^\perp}, dy_{F^\perp})) \gamma_{E^\perp \times F^\perp | E \times F}^*((x'_E, y'_F), (dx'_{E^\perp}, dy'_{F^\perp})) \end{aligned}$$

by definition of the Monge-Knothe coupling. This is well optimal for subspace optimal plans. \square

Proof of Proposition 2. Let $f : \mathbb{R}^p \rightarrow \mathbb{R}^p$ be an invariance of GW on E^\perp , i.e. $\forall x \in \mathbb{R}^p, f(x) = (x_E, f_{E^\perp}(x_{E^\perp}))$. We first deal with $L(x, x', y, y') = (\|x - x'\|_2^2 - \|y - y'\|_2^2)^2$, and therefore f_{E^\perp} is either an isometry or a translation.

From lemma 6 of Paty and Cuturi (2019), we know that $\Pi(f\#\mu, \nu) = \{(f, Id)\#\gamma \mid \gamma \in \Pi(\mu, \nu)\}$. We can rewrite

$$\begin{aligned} \Pi_{E,F}(f\#\mu, \nu) &= \{\gamma \in \Pi(f\#\mu, \nu) \mid (\pi^E, \pi^F)\#\gamma = \gamma_{E,F}^*\} \\ &= \{(f, Id)\#\gamma \mid \gamma \in \Pi(\mu, \nu), (\pi^E, \pi^F)\#\gamma = \gamma_{E,F}^*\} \\ &= \{(f, Id)\#\gamma \mid \gamma \in \Pi(\mu, \nu), (\pi^E, \pi^F)\#\gamma = \gamma_{E,F}^*\} \\ &= \{(f, Id)\#\gamma \mid \gamma \in \Pi_{E,F}(\mu, \nu)\} \end{aligned}$$

using $f = (Id_E, f_{E^\perp})$, $\pi^E \circ f = Id_E$ and $(\pi^E, \pi^F)\#\gamma = (\pi^E, \pi^F)\#\gamma$.

Now, for all $\gamma \in \Pi_{E,F}(f\#\mu, \nu)$, there exists $\tilde{\gamma} \in \Pi_{E,F}(\mu, \nu)$ such that $\gamma = (f, Id)\#\tilde{\gamma}$ and we can disintegrate $\tilde{\gamma}$ with respect to $\gamma_{E,F}^*$

$$\tilde{\gamma} = \gamma_{E,F}^* \otimes K$$

with K a probability kernel on $(E \times F, \mathcal{B}(E) \otimes \mathcal{B}(F))$.

For $\gamma_{E,F}^*$ almost every $(x_E, y_F), (x'_E, y'_F)$, we have

$$\begin{aligned} & \iint (\|x_E - x'_E\|_2^2 + \|x_{E^\perp} - x'_{E^\perp}\|_2^2 - \|y_F - y'_F\|_2^2 - \|y_{F^\perp} - y'_{F^\perp}\|_2^2)^2 \\ & (f_{E^\perp}, Id)\#K((x_E, y_F), (dx_{E^\perp}, dy_{F^\perp})) (f_{E^\perp}, Id)\#K((x'_E, y'_F), (dx'_{E^\perp}, dy'_{F^\perp})) \\ &= \iint (\|x_E - x'_E\|_2^2 + \|f_{E^\perp}(x_{E^\perp}) - f_{E^\perp}(x'_{E^\perp})\|_2^2 - \|y_F - y'_F\|_2^2 - \|y_{F^\perp} - y'_{F^\perp}\|_2^2)^2 \\ & K((x_E, y_F), (dx_{E^\perp}, dy_{F^\perp})) K((x'_E, y'_F), (dx'_{E^\perp}, dy'_{F^\perp})) \\ &= \iint (\|x_E - x'_E\|_2^2 + \|x_{E^\perp} - x'_{E^\perp}\|_2^2 - \|y_F - y'_F\|_2^2 - \|y_{F^\perp} - y'_{F^\perp}\|_2^2)^2 \\ & K((x_E, y_F), (dx_{E^\perp}, dy_{F^\perp})) K((x'_E, y'_F), (dx'_{E^\perp}, dy'_{F^\perp})) \end{aligned}$$

using in the last line that $\|f_{E^\perp}(x_{E^\perp}) - f_{E^\perp}(x'_{E^\perp})\|_2 = \|x_{E^\perp} - x'_{E^\perp}\|_2$ since $d(x, y) = \|x - y\|_2$ is translation and rotation invariant ($d(Ox, Oy) = d(x, y)$ and $d(Tx, Ty) = d(x, y)$).

By integrating with respect to $\gamma_{E,F}^*$, we obtain

$$\begin{aligned}
& \iint \left(\iint (\|x - x'\|_2^2 - \|y - y'\|_2^2)^2 \right. \\
& \left. (f_{E^\perp}, Id)_{\#} K((x_E, y_F), (dx_{E^\perp}, dy_{F^\perp})) (f_{E^\perp}, Id)_{\#} K((x'_E, y'_F), (dx'_{E^\perp}, dy'_{F^\perp})) \right) d\gamma_{E,F}^*(x_E, y_F) d\gamma_{E,F}^*(x'_E, y'_F) \\
& = \iint (\|x - x'\|_2^2 - \|y - y'\|_2^2)^2 d\tilde{\gamma}(x, y) d\tilde{\gamma}(x', y').
\end{aligned} \tag{9}$$

Now, we show that $\gamma = (f, Id)_{\#} \tilde{\gamma} = \gamma_{E,F}^* \otimes (f_{E^\perp}, Id)_{\#} K$. Let ϕ some bounded measurable function on $\mathbb{R}^p \times \mathbb{R}^q$,

$$\begin{aligned}
\int \phi(x, y) d\gamma(x, y) &= \int \phi(x, y) d((f, Id)_{\#} \tilde{\gamma}(x, y)) \\
&= \int \phi(f(x), y) d\tilde{\gamma}(x, y) \\
&= \iint \phi(f(x), y) K((x_E, y_F), (dx_{E^\perp}, dy_{F^\perp})) d\gamma_{E,F}^*(x_E, y_F) \\
&= \iint \phi((x_E, f_{E^\perp}(x_{E^\perp})), y) K((x_E, y_F), (dx_{E^\perp}, dy_{F^\perp})) d\gamma_{E,F}^*(x_E, y_F) \\
&= \iint \phi(x, y) (f_{E^\perp}, Id)_{\#} K((x_E, y_F), (dx_{E^\perp}, dy_{F^\perp})) d\gamma_{E,F}^*(x_E, y_F).
\end{aligned}$$

Hence, we can rewrite (9) as

$$\iint (\|x - x'\|_2^2 - \|y - y'\|_2^2)^2 d(f, Id)_{\#} \tilde{\gamma}(x, y) d(f, Id)_{\#} \tilde{\gamma}(x', y') = \iint (\|x - x'\|_2^2 - \|y - y'\|_2^2)^2 d\tilde{\gamma}(x, y) d\tilde{\gamma}(x', y').$$

Now, by taking the infimum with respect to $\tilde{\gamma} \in \Pi_{E,F}(\mu, \nu)$, we find

$$GW_{E,F}^2(f_{\#}\mu, \nu) = GW_{E,F}^2(\mu, \nu).$$

For the inner product case, we can do the same proof for isometries. \square

A.2 Closed-form between Gaussians

Let $q \leq p$, $\mu = \mathcal{N}(m_\mu, \Sigma) \in \mathcal{P}(\mathbb{R}^p)$, $\nu = \mathcal{N}(m_\nu, \Lambda) \in \mathcal{P}(\mathbb{R}^q)$ two Gaussian measures with $\Sigma = P_\mu D_\mu P_\mu^T$ and $\Lambda = P_\nu D_\nu P_\nu^T$.

Let $E \subset \mathbb{R}^p$ be a subspace of dimension k and $F \subset \mathbb{R}^q$ a subspace of dimension k' .

We represent Σ in an orthonormal basis of $E \oplus E^\perp$, and Λ in an orthonormal basis of $F \oplus F^\perp$, *i.e.*

$\Sigma = \begin{pmatrix} \Sigma_E & \Sigma_{EE^\perp} \\ \Sigma_{E^\perp E} & \Sigma_{E^\perp E^\perp} \end{pmatrix}$. We denote $\Sigma/\Sigma_E = \Sigma_{E^\perp} - \Sigma_{E^\perp E}^T \Sigma_E^{-1} \Sigma_{EE^\perp}$ the Schur complement of Σ with respect to Σ_E . We know that the conditionals of Gaussians are Gaussians, and of covariance the Schur complement (see *e.g.* Rasmussen (2003) or Von Mises (1964)).

A.2.1 Quadratic GW problem

For GW with $c(x, x') = \|x - x'\|_2^2$, we have for now no guarantee that there exists an optimal coupling which is a transport map. Salmona et al. proposed to restrict the problem to the set of Gaussian couplings $\pi(\mu, \nu) \cap \mathcal{N}_{p+q}$ where \mathcal{N}_{p+q} denotes the set of Gaussians in \mathbb{R}^{p+q} . In that case, the problem becomes

$$GGW_2^2(\mu, \nu) = \inf_{\gamma \in \Pi(\mu, \nu) \cap \mathcal{N}_{p+q}} \iint (\|x - x'\|_2^2 - \|y - y'\|_2^2)^2 d\gamma(x, y) d\gamma(x', y'). \tag{10}$$

In that case, they showed that an optimal solution is of the form $T(x) = m_\nu + P_\nu A P_\mu^T (x - m_\mu)$ with $A = \left(\tilde{I}_q D_\nu^{\frac{1}{2}} (D_\mu^{(q)})^{-\frac{1}{2}} \quad 0_{q, p-q} \right)$ and \tilde{I}_q of the form $\text{diag}((\pm 1)_{i \leq q})$.

Since the problem is translation invariant, we can always solve the problem between the centered measures.

In the following, we suppose that $k = k'$. Let's denote $T_{E,F}$ the optimal transport map for (10) between $\mathcal{N}(0, \Sigma_E)$ and $\mathcal{N}(0, \Lambda_F)$. According to Theorem 4.1 in Salmona et al. (2021), such a solution exists and is of the form (5). We also denote T_{E^\perp, F^\perp} the optimal transport map between $\mathcal{N}(0, \Sigma/\Sigma_E)$ and $\mathcal{N}(0, \Lambda/\Lambda_F)$ (which is well defined since we assumed $p \geq q$ and hence $p - k \geq q - k'$ since $k = k'$).

We know that the Monge-Knothe transport map will be a linear map $T_{\text{MK}}(x) = Bx$ with B a block triangular matrix of the form

$$B = \begin{pmatrix} T_{E,F} & 0_{k', p-k} \\ C & T_{E^\perp, F^\perp} \end{pmatrix} \in \mathbb{R}^{q \times p},$$

with $C \in \mathbb{R}^{(q-k') \times k}$, and such that $B\Sigma B^T = \Lambda$ (to have well a transport map between μ and ν).

Actually,

$$B\Sigma B^T = \begin{pmatrix} T_{E,F}\Sigma_E T_{E,F}^T & T_{E,F}\Sigma_E C^T + T_{E,F}\Sigma_{EE^\perp} T_{E^\perp, F^\perp}^T \\ ((C\Sigma_E + T_{E^\perp, F^\perp}\Sigma_{E^\perp E})T_{E,F}^T & (C\Sigma_E + T_{E^\perp, F^\perp}\Sigma_{E^\perp E})C^T + (C\Sigma_{EE^\perp} + T_{E^\perp, F^\perp}\Sigma_{E^\perp})T_{E^\perp, F^\perp}^T \end{pmatrix}$$

First, we have well $T_{E,F}\Sigma_E T_{E,F}^T = \Lambda_F$ as $T_{E,F}$ is a transport map between μ_E and ν_F . Then,

$$B\Sigma B^T = \Lambda \iff \begin{cases} T_{E,F}\Sigma_E T_{E,F}^T = \Lambda_F \\ T_{E,F}\Sigma_E C^T + T_{E,F}\Sigma_{EE^\perp} T_{E^\perp, F^\perp}^T = \Lambda_{FF^\perp} \\ (C\Sigma_E + T_{E^\perp, F^\perp}\Sigma_{E^\perp E})T_{E,F}^T = \Lambda_{F^\perp F} \\ (C\Sigma_E + T_{E^\perp, F^\perp}\Sigma_{E^\perp E})C^T + (C\Sigma_{EE^\perp} + T_{E^\perp, F^\perp}\Sigma_{E^\perp})T_{E^\perp, F^\perp}^T = \Lambda_{F^\perp}. \end{cases}$$

We have

$$(C\Sigma_E + T_{E^\perp, F^\perp}\Sigma_{E^\perp E})T_{E,F}^T = \Lambda_{F^\perp F} \iff C\Sigma_E T_{E,F}^T = \Lambda_{F^\perp F} - T_{E^\perp, F^\perp}\Sigma_{E^\perp E} T_{E,F}^T.$$

As $k = k'$, $\Sigma_E T_{E,F}^T \in \mathbb{R}^{k \times k}$ and is invertible (as Σ_E and Λ_F are positive definite and $T_{E,F} = P_{\mu_E} A_{E,F} P_{\nu_F}$ with $A_{E,F} = \left(\tilde{I}_k D_{\nu_F}^{\frac{1}{2}} D_{\mu_E}^{-\frac{1}{2}} \right)$ with positive values on the diagonals. Hence, we have

$$C = (\Lambda_{F^\perp F} (T_{E,F}^T)^{-1} - T_{E^\perp, F^\perp}\Sigma_{E^\perp E})\Sigma_E^{-1}.$$

Now, we still have to check the last two equations. First,

$$\begin{aligned} T_{E,F}\Sigma_E C^T + T_{E,F}\Sigma_{EE^\perp} T_{E^\perp, F^\perp}^T &= T_{E,F}\Sigma_E \Sigma_E^{-1} T_{E,F}^{-1} \Lambda_{F^\perp F}^T - T_{E,F}\Sigma_E \Sigma_E^{-1} \Sigma_{E^\perp E} T_{E^\perp, F^\perp}^T + T_{E,F}\Sigma_{EE^\perp} T_{E^\perp, F^\perp}^T \\ &= \Lambda_{FF^\perp}. \end{aligned}$$

And for the last equation,

$$\begin{aligned} &(C\Sigma_E + T_{E^\perp, F^\perp}\Sigma_{E^\perp E})C^T + (C\Sigma_{EE^\perp} + T_{E^\perp, F^\perp}\Sigma_{E^\perp})T_{E^\perp, F^\perp}^T \\ &= (\Lambda_{F^\perp F} (T_{E,F}^T)^{-1} - T_{E^\perp, F^\perp}\Sigma_{E^\perp E} + T_{E^\perp, F^\perp}\Sigma_{E^\perp E})\Sigma_E^{-1} (T_{E,F}^{-1} \Lambda_{F^\perp F}^T - \Sigma_{E^\perp E} T_{E^\perp, F^\perp}^T) \\ &\quad + \Lambda_{F^\perp F} (T_{E,F}^T)^{-1} \Sigma_E^{-1} \Sigma_{EE^\perp} T_{E^\perp, F^\perp}^T - T_{E^\perp, F^\perp}\Sigma_{E^\perp E} \Sigma_E^{-1} \Sigma_{EE^\perp} T_{E^\perp, F^\perp}^T + T_{E^\perp, F^\perp}\Sigma_{E^\perp} T_{E^\perp, F^\perp}^T \\ &= \Lambda_{F^\perp F} (T_{E,F}^T)^{-1} \Sigma_E^{-1} T_{E,F}^{-1} \Lambda_{F^\perp F}^T - \Lambda_{F^\perp F} (T_{E,F}^T)^{-1} \Sigma_E^{-1} \Sigma_{E^\perp E} T_{E^\perp, F^\perp}^T - T_{E^\perp, F^\perp}\Sigma_{E^\perp E} \Sigma_E^{-1} T_{E,F}^{-1} \Lambda_{F^\perp F}^T \\ &\quad + T_{E^\perp, F^\perp}\Sigma_{E^\perp E} \Sigma_E^{-1} \Sigma_{E^\perp E} T_{E^\perp, F^\perp}^T + T_{E^\perp, F^\perp}\Sigma_{E^\perp E} \Sigma_E^{-1} T_{E,F}^{-1} \Lambda_{F^\perp F}^T - T_{E^\perp, F^\perp}\Sigma_{E^\perp E} \Sigma_E^{-1} \Sigma_{E^\perp E} T_{E^\perp, F^\perp}^T \\ &\quad + \Lambda_{F^\perp F} (T_{E,F}^T)^{-1} \Sigma_E^{-1} \Sigma_{EE^\perp} T_{E^\perp, F^\perp}^T - T_{E^\perp, F^\perp}\Sigma_{E^\perp E} \Sigma_E^{-1} \Sigma_{E^\perp E} T_{E^\perp, F^\perp}^T + T_{E^\perp, F^\perp}\Sigma_{E^\perp} T_{E^\perp, F^\perp}^T \\ &= \Lambda_{F^\perp F} (T_{E,F}^T)^{-1} \Sigma_E^{-1} T_{E,F}^{-1} \Lambda_{F^\perp F}^T - T_{E^\perp, F^\perp}\Sigma_{E^\perp E} \Sigma_E^{-1} \Sigma_{E^\perp E} T_{E^\perp, F^\perp}^T + T_{E^\perp, F^\perp}\Sigma_{E^\perp} T_{E^\perp, F^\perp}^T \end{aligned}$$

Now, using that $(T_{E,F}^T)^{-1} \Sigma_E^{-1} T_{E,F}^{-1} = (T_{E,F}\Sigma_E T_{E,F}^T)^{-1} = \Lambda_F^{-1}$ and $\Sigma_{E^\perp} - \Sigma_{E^\perp E} \Sigma_E^{-1} \Sigma_{E^\perp E} = \Sigma/\Sigma_E$, we have

$$\begin{aligned} &(C\Sigma_E + T_{E^\perp, F^\perp}\Sigma_{E^\perp E})C^T + (C\Sigma_{EE^\perp} + T_{E^\perp, F^\perp}\Sigma_{E^\perp})T_{E^\perp, F^\perp}^T \\ &= \Lambda_{F^\perp F} \Lambda_F^{-1} \Lambda_{F^\perp F}^T + T_{E^\perp, F^\perp} (\Sigma_{E^\perp} - \Sigma_{E^\perp E} \Sigma_E^{-1} \Sigma_{E^\perp E}) T_{E^\perp, F^\perp}^T \\ &= \Lambda_{F^\perp F} \Lambda_F^{-1} \Lambda_{F^\perp F}^T + \Lambda/\Lambda_F \\ &= \Lambda_{F^\perp} \end{aligned}$$

Then π_{MK} is of the form $(Id, T_{\text{MK}})_{\#}\mu$ with

$$T_{\text{MK}}(x) = m_{\nu} + B(x - m_{\mu}).$$

A.2.2 Closed-form between Gaussians for Monge-Independent

Suppose $k \geq k'$ in order to be able to define the OT map between μ_E and ν_F .

For the Monge-Independent plan $\pi_{\text{MI}} = \gamma_{E,F}^* \otimes (\mu_{E^{\perp}|E} \otimes \nu_{F^{\perp}|F})$, let $(X, Y) \sim \pi_{\text{MI}}$. We know that π_{MI} is a degenerate Gaussian with a covariance of the form

$$\text{Cov}(X, Y) = \begin{pmatrix} \text{Cov}(X) & C \\ C^T & \text{Cov}(Y) \end{pmatrix}$$

where $\text{Cov}(X) = \Sigma$ and $\text{Cov}(Y) = \Lambda$. Moreover, we know that C is of the form

$$\begin{pmatrix} \text{Cov}(X_E, Y_F) & \text{Cov}(X_E, Y_{F^{\perp}}) \\ \text{Cov}(X_{E^{\perp}}, Y_F) & \text{Cov}(X_{E^{\perp}}, Y_{F^{\perp}}) \end{pmatrix}.$$

Let's assume that $m_{\mu} = m_{\nu} = 0$, then

$$\text{Cov}(X_E, Y_F) = \text{Cov}(X_E, T_{E,F}X_E) = \mathbb{E}[X_E X_E^T] T_{E,F}^T = \Sigma_E T_{E,F}^T,$$

$$\begin{aligned} \text{Cov}(X_E, Y_{F^{\perp}}) &= \mathbb{E}[X_E Y_{F^{\perp}}^T] \\ &= \mathbb{E}[\mathbb{E}[X_E Y_{F^{\perp}}^T | X_E, Y_F]] \\ &= \mathbb{E}[X_E \mathbb{E}[Y_{F^{\perp}}^T | Y_F]] \end{aligned}$$

since $Y_F = T_{E,F}X_E$, X_E is $\sigma(Y_F)$ -measurable. Now, using the equation (A.6) from Rasmussen (2003), we have

$$\begin{aligned} \mathbb{E}[Y_{F^{\perp}} | Y_F] &= \Lambda_{F^{\perp}F} \Lambda_F^{-1} Y_F \\ &= \Lambda_{F^{\perp}F} \Lambda_F^{-1} T_{E,F} X_E \end{aligned}$$

and

$$\mathbb{E}[X_{E^{\perp}} | X_E] = \Sigma_{E^{\perp}E} \Sigma_E^{-1} X_E.$$

Hence,

$$\begin{aligned} \text{Cov}(X_E, Y_{F^{\perp}}) &= \mathbb{E}[X_E \mathbb{E}[Y_{F^{\perp}}^T | Y_F]] \\ &= \mathbb{E}[X_E X_E^T] T_{E,F}^T \Lambda_F^{-1} \Lambda_{F^{\perp}F}^T \\ &= \Sigma_E T_{E,F}^T \Lambda_F^{-1} \Lambda_{F^{\perp}F}^T. \end{aligned}$$

We also have

$$\text{Cov}(X_{E^{\perp}}, Y_F) = \mathbb{E}[X_{E^{\perp}} X_E^T T_{E,F}^T] = \Sigma_{E^{\perp}E} T_{E,F}^T,$$

and

$$\begin{aligned} \text{Cov}(X_{E^{\perp}}, Y_{F^{\perp}}) &= \mathbb{E}[X_{E^{\perp}} Y_{F^{\perp}}^T] \\ &= \mathbb{E}[\mathbb{E}[X_{E^{\perp}} Y_{F^{\perp}}^T | X_E, Y_F]] \\ &= \mathbb{E}[\mathbb{E}[X_{E^{\perp}} | X_E] \mathbb{E}[Y_{F^{\perp}}^T | Y_F]] \text{ by independence} \\ &= \mathbb{E}[\Sigma_{E^{\perp}E} \Sigma_E^{-1} X_E X_E^T T_{E,F}^T \Lambda_F^{-1} \Lambda_{F^{\perp}F}^T] \\ &= \Sigma_{E^{\perp}E} T_{E,F}^T \Lambda_F^{-1} \Lambda_{F^{\perp}F}^T. \end{aligned}$$

Finally, we find

$$C = \begin{pmatrix} \Sigma_E T_{E,F}^T & \Sigma_E T_{E,F}^T \Lambda_F^{-1} \Lambda_{F^{\perp}F}^T \\ \Sigma_{E^{\perp}E} T_{E,F}^T & \Sigma_{E^{\perp}E} T_{E,F}^T \Lambda_F^{-1} \Lambda_{F^{\perp}F}^T \end{pmatrix}.$$

By taking orthogonal bases $(V_E, V_{E^{\perp}})$ and $(V_F, V_{F^{\perp}})$, we can put it in a more compact way such as in Proposition 4 in Muzellec and Cuturi (2019):

$$C = (V_E \Sigma_E + V_{E^{\perp}} \Sigma_{E^{\perp}E}) T_{E,F}^T (V_F^T + \Lambda_F^{-1} \Lambda_{F^{\perp}F}^T V_{F^{\perp}}^T).$$

To check it, just expand the terms and see that $C_{E,F} = V_E C V_F^T$.

B Knothe-Rosenblatt

B.1 Properties of (7)

Proof of Proposition 5. Let $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$,

1. $(x, x') \mapsto x \odot x'$ is a continuous map, therefore L is lower semi-continuous. Hence, by applying lemma 2.2.1 of (Vayer, 2020), we have that $\gamma \mapsto \iint L(x, x', y, y') d\gamma(x, y) d\gamma(x', y')$ is lower semi-continuous for the weak convergence of measures.

Now, as $\Pi(\mu, \nu)$ is a compact set (see the proof of Theorem 1.7 in Santambrogio (2015) for the Polish space case, and of Theorem 1.4 for the compact metric space), and $\gamma \mapsto \iint L d\gamma d\gamma$ is lower semi-continuous for the weak convergence, we can apply the Weierstrass theorem (Memo 2.2.1 in Vayer (2020)) which states that (7) always admits a minimizer.

2. See Theorem 16 in Chowdhury and Mémoli (2019).
3. For invariances, we first look at the properties that must be satisfied by T in order to have: $\forall x, x', f(x, x') = f(T(x), T(x'))$ where $f : (x, x') \mapsto x \odot x'$.

We find that $\forall x \in \mathbb{R}^d, \forall 1 \leq i \leq d, |[T(x)]_i| = |x_i|$ because, denoting $(e_i)_{i=1}^d$ as the canonical basis, we have

$$f(e_i, x) = x e_i = f(T(e_i), T(x)) \implies [T(e_i)]_i [T(x)]_i = x_i \implies |[T(x)]_i| = |x_i|$$

$$\text{as } f(e_i, e_i) = f(T(e_i), T(e_i)) \implies [T(e_i)]_i^2 = 1.$$

If we take for T the reflection with respect to axis, then it satisfies well $f(x, x') = f(T(x), T(x'))$. Moreover, it is well an equivalent relation, and therefore we have a distance on the quotient space.

□

Proposition 6. *In a slightly more general setting, let $X_0 = X_1 = \mathbb{R}^d$, functions f_0, f_1 from $\mathbb{R}^d \times \mathbb{R}^d$ to \mathbb{R}^d and measures $\mu_0 \in \mathcal{P}(X_0), \mu_1 \in \mathcal{P}(X_1)$. Then the family $\mathcal{X}_t = (X_0 \times X_1, f_t, \gamma^*)$ defines a geodesic between \mathcal{X}_0 and \mathcal{X}_1 , where γ^* is the optimal coupling of \mathcal{HW} between μ_0 and μ_1 , and*

$$f_t((x_0, x'_0), (x_1, x'_1)) = (1 - t)f_0(x_0, x'_0) + t f_1(x_1, x'_1).$$

Proof. See Theorem 3.1 in Sturm (2012). □

B.2 Proof of Theorem 3

We first recall a useful theorem.

Theorem 4 (Theorem 2.8 in Billingsley (2013)). *Let $\Omega = X \times Y$ be a separable space, and let $P, P_n \in \mathcal{P}(\Omega)$ with marginals P_X (respectively $P_{n,X}$) and P_Y (respectively $P_{n,Y}$). Then, $P_{n,X} \otimes P_{n,Y} \xrightarrow{\mathcal{D}} P$ if and only if $P_{n,X} \xrightarrow{\mathcal{D}} P_X, P_{n,Y} \xrightarrow{\mathcal{D}} P_Y$ and $P = P_X \otimes P_Y$.*

Proof of Theorem 3. The following proof is mainly inspired by the proof of Theorem 1 in (Carlier et al., 2010)[Theorem 2.1], (Bonnotte, 2013)[Theorem 3.1.6] and (Santambrogio, 2015)[Theorem 2.23].

Let $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$, absolutely continuous, with finite fourth moments and compact supports. We recall the problem \mathcal{HW}_t ,

$$\mathcal{HW}_t^2(\mu, \nu) = \inf_{\gamma \in \Pi(\mu, \nu)} \iint \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_t(x, y) d\gamma_t(x', y'),$$

with $\forall t > 0, \forall i \in \{1, \dots, d-1\}, \lambda_t^{(i)} > 0$ and $\lambda_t^{(i)} \xrightarrow{t \rightarrow 0} 0$.

First, let's denote γ_t the optimal coupling for \mathcal{HW}_t for all $t > 0$. We want to show that $\gamma_t \xrightarrow[t \rightarrow 0]{\mathcal{D}} \gamma_K$ with $\gamma_K = (Id \times T_K) \# \mu$ and T_K our alternate Knothe-Rosenblatt rearrangement. Let $\gamma \in \Pi(\mu, \nu)$

such that $\gamma_t \xrightarrow[t \rightarrow 0]{\mathcal{D}} \gamma$ (true up to subsequence as $\{\mu\}$ and $\{\nu\}$ are tight in $\mathcal{P}(X)$ and $\mathcal{P}(Y)$ if X and Y are polish space, therefore, by (Villani, 2008)[Lemma 4.4], $\Pi(\mu, \nu)$ is a tight set, and we can apply the Prokhorov theorem (Santambrogio, 2015)[Box 1.4] on $(\gamma_t)_t$ and extract a subsequence).

Part 1:

First, let's notice that:

$$\begin{aligned} \mathcal{HW}_t^2(\mu, \nu) &= \iint \sum_{k=1}^d \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_t(x, y) d\gamma_t(x', y') \\ &= \iint (x_1 x'_1 - y_1 y'_1)^2 d\gamma_t(x, y) d\gamma_t(x', y') + \iint \sum_{k=2}^d \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_t(x, y) d\gamma_t(x', y'). \end{aligned}$$

Moreover, as γ_t is the optimal coupling between μ and ν , and $\gamma_K \in \Pi(\mu, \nu)$,

$$\begin{aligned} \mathcal{HW}_t^2(\mu, \nu) &\leq \iint \sum_{k=1}^d \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_K(x, y) d\gamma_K(x', y') \\ &= \iint (x_1 x'_1 - y_1 y'_1)^2 d\gamma_K(x, y) d\gamma_K(x', y') + \iint \sum_{k=2}^d \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_K(x, y) d\gamma_K(x', y'). \end{aligned}$$

In our case, we have $\gamma_t \xrightarrow[t \rightarrow 0]{\mathcal{D}} \gamma$, thus, by Theorem 4, we have $\gamma_t \otimes \gamma_t \xrightarrow[t \rightarrow 0]{\mathcal{D}} \gamma \otimes \gamma$. Using the fact that $\forall i, \lambda_t^{(i)} \xrightarrow[t \rightarrow 0]{} 0$ (and lemma 1.8 of Santambrogio (2015), since we are on compact support, we can bound the cost (which is continuous) by its max), we obtain the following inequality

$$\iint (x_1 x'_1 - y_1 y'_1)^2 d\gamma(x, y) d\gamma(x', y') \leq \iint (x_1 x'_1 - y_1 y'_1)^2 d\gamma_K(x, y) d\gamma_K(x', y').$$

By denoting γ^1 and γ_K^1 the marginals on the first variables, we can use the projection $\pi^1(x, y) = (x_1, y_1)$, such as $\gamma^1 = \pi_{\#}^1 \gamma$ and $\gamma_K^1 = \pi_{\#}^1 \gamma_K$. Hence, we get

$$\iint (x_1 x'_1 - y_1 y'_1)^2 d\gamma^1(x_1, y_1) d\gamma^1(x'_1, y'_1) \leq \iint (x_1 x'_1 - y_1 y'_1)^2 d\gamma_K^1(x_1, y_1) d\gamma_K^1(x'_1, y'_1).$$

However, γ_K^1 was constructed in order to be the unique optimal map for this cost (either T_{asc} or T_{desc} according to theorem (Vayer, 2020)[Theorem 4.2.4]). Thus, we can deduce that $\gamma^1 = (Id \times T_K^1)_{\#} \mu^1 = \gamma_K^1$.

Part 2:

We know that for any $t > 0$, γ_t and γ_K share the same marginals. Thus, as previously, $\pi_{\#}^1 \gamma_t$ should have a cost worse than $\pi_{\#}^1 \gamma_K$, which translates to

$$\begin{aligned} \iint (x_1 x'_1 - y_1 y'_1)^2 d\gamma_K^1(x_1, y_1) d\gamma_K^1(x'_1, y'_1) &= \iint (x_1 x'_1 - y_1 y'_1)^2 d\gamma^1(x_1, y_1) d\gamma^1(x'_1, y'_1) \\ &\leq \iint (x_1 x'_1 - y_1 y'_1)^2 d\gamma_t^1(x_1, y_1) d\gamma_t^1(x'_1, y'_1). \end{aligned}$$

Therefore, we have the following inequality,

$$\begin{aligned} &\iint (x_1 x'_1 - y_1 y'_1)^2 d\gamma^1(x, y) d\gamma^1(x', y') + \iint \sum_{k=2}^d \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_t(x, y) d\gamma_t(x', y') \\ &\leq \mathcal{HW}_t^2(\mu, \nu) \\ &\leq \iint (x_1 x'_1 - y_1 y'_1)^2 d\gamma^1(x, y) d\gamma^1(x', y') + \iint \sum_{k=2}^d \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_K(x, y) d\gamma_K(x', y'). \end{aligned}$$

We can subtract the first term and factorize by $\lambda_t^{(1)} > 0$,

$$\begin{aligned}
& \iint \sum_{k=2}^d \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_t(x, y) d\gamma_t(x', y') \\
&= \lambda_t^{(1)} \left(\iint (x_2 x'_2 - y_2 y'_2)^2 d\gamma_t(x, y) d\gamma_t(x', y') + \iint \sum_{k=3}^d \left(\prod_{i=2}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_t(x, y) d\gamma_t(x', y') \right) \\
&\leq \lambda_t^{(1)} \left(\iint (x_2 x'_2 - y_2 y'_2)^2 d\gamma_K(x, y) d\gamma_K(x', y') + \iint \sum_{k=3}^d \left(\prod_{i=2}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_K(x, y) d\gamma_K(x', y') \right).
\end{aligned}$$

By dividing by $\lambda_t^{(1)}$ and by taking the limit $t \rightarrow 0$ as in the first part, we get

$$\iint (x_2 x'_2 - y_2 y'_2)^2 d\gamma(x, y) d\gamma(x', y') \leq \iint (x_2 x'_2 - y_2 y'_2)^2 d\gamma_K(x, y) d\gamma_K(x', y'). \quad (11)$$

Now, the 2 terms depend only on (x_2, y_2) and (x'_2, y'_2) . We will project on the two first coordinates, *i.e.* let $\pi^{1,2}(x, y) = ((x_1, x_2), (y_1, y_2))$ and $\gamma^{1,2} = \pi_{\#}^{1,2} \gamma$, $\gamma_K^{1,2} = \pi_{\#}^{1,2} \gamma_K$. Using the disintegration of measures, we know that there exist kernels $\gamma^{2|1}$ and $\gamma_K^{2|1}$ such that $\gamma^{1,2} = \gamma^1 \otimes \gamma^{2|1}$ and $\gamma_K^{1,2} = \gamma_K^1 \otimes \gamma_K^{2|1}$, where

$$\forall A \in \mathcal{B}(X \times Y), \mu \otimes K(A) = \iint \mathbb{1}_A(x, y) K(x, dy) \mu(dx).$$

We can rewrite the previous equation (11) as

$$\begin{aligned}
& \iint (x_2 x'_2 - y_2 y'_2)^2 d\gamma(x, y) d\gamma(x', y') \\
&= \iiint \iint (x_2 x'_2 - y_2 y'_2)^2 \gamma^{2|1}((x_1, y_1), (dx_2, dy_2)) \gamma^{2|1}((x'_1, y'_1), (dx'_2, dy'_2)) d\gamma^1(x_1, y_1) d\gamma^1(x'_1, y'_1) \\
&\leq \iiint \iint (x_2 x'_2 - y_2 y'_2)^2 \gamma_K^{2|1}((x_1, y_1), (dx_2, dy_2)) \gamma_K^{2|1}((x'_1, y'_1), (dx'_2, dy'_2)) d\gamma_K^1(x_1, y_1) d\gamma_K^1(x'_1, y'_1).
\end{aligned} \quad (12)$$

Now, we will assume at first that the marginals of $\gamma^{2|1}((x_1, y_1), \cdot)$ are well $\mu^{2|1}(x_1, \cdot)$ and $\nu^{2|1}(y_1, \cdot)$. Then, by definition of $\gamma_K^{2|1}$, as it is optimal for the *GW* cost with inner products, we have for all $(x_1, y_1), (x'_1, y'_1)$,

$$\begin{aligned}
& \iint (x_2 x'_2 - y_2 y'_2)^2 \gamma_K^{2|1}((x_1, y_1), (dx_2, dy_2)) \gamma_K^{2|1}((x'_1, y'_1), (dx'_2, dy'_2)) \\
&\leq \iint (x_2 x'_2 - y_2 y'_2)^2 \gamma^{2|1}((x_1, y_1), (dx_2, dy_2)) \gamma^{2|1}((x'_1, y'_1), (dx'_2, dy'_2)).
\end{aligned} \quad (13)$$

Moreover, we know from the first part that $\gamma^1 = \gamma_K^1$, then by integrating with respect to (x_1, y_1) and (x'_1, y'_1) , we have

$$\begin{aligned}
& \iiint \iint (x_2 x'_2 - y_2 y'_2)^2 \gamma_K^{2|1}((x_1, y_1), (dx_2, dy_2)) \gamma_K^{2|1}((x'_1, y'_1), (dx'_2, dy'_2)) d\gamma^1(x_1, y_1) d\gamma^1(x'_1, y'_1) \\
&\leq \iiint \iint (x_2 x'_2 - y_2 y'_2)^2 \gamma^{2|1}((x_1, y_1), (dx_2, dy_2)) \gamma^{2|1}((x'_1, y'_1), (dx'_2, dy'_2)) d\gamma^1(x_1, y_1) d\gamma^1(x'_1, y'_1).
\end{aligned} \quad (14)$$

By (12) and (14), we deduce that we have an equality and we get

$$\begin{aligned}
& \iint \left(\iint (x_2 x'_2 - y_2 y'_2)^2 \gamma^{2|1}((x_1, y_1), (dx_2, dy_2)) \gamma^{2|1}((x'_1, y'_1), (dx'_2, dy'_2)) \right. \\
& \left. - \iint (x_2 x'_2 - y_2 y'_2)^2 \gamma_K^{2|1}((x_1, y_1), (dx_2, dy_2)) \gamma_K^{2|1}((x'_1, y'_1), (dx'_2, dy'_2)) \right) d\gamma^1(x_1, y_1) d\gamma^1(x'_1, y'_1) = 0.
\end{aligned} \quad (15)$$

However, we know by (13) that the middle part of (15) is nonnegative, thus we have for γ^1 -a.e. $(x_1, y_1), (x'_1, y'_1)$,

$$\begin{aligned} & \iint (x_2 x'_2 - y_2 y'_2)^2 \gamma_K^{2|1}((x_1, y_1), (dx_2, dy_2)) \gamma_K^{2|1}((x'_1, y'_1), (dx'_2, dy'_2)) \\ &= \iint (x_2 x'_2 - y_2 y'_2)^2 \gamma^{2|1}((x_1, y_1), (dx_2, dy_2)) \gamma^{2|1}((x'_1, y'_1), (dx'_2, dy'_2)). \end{aligned}$$

From that, we can conclude as in the first part that $\gamma^{2|1} = \gamma_K^{2|1}$ (by unicity of the optimal map). And thus $\gamma^{1,2} = \gamma_K^{1,2}$.

Now, we still have to show that the marginals of $\gamma^{2|1}((x_1, y_1), \cdot)$ and $\gamma_K^{2|1}((x_1, y_1), \cdot)$ are well the same, i.e. $\mu^{2|1}(x_1, \cdot)$ and $\nu^{2|1}(y_1, \cdot)$. Let ϕ and ψ be continuous functions, then we have to show that for γ^1 -a.e. (x_1, y_1) , we have

$$\begin{cases} \int \phi(x_2) \gamma^{2|1}((x_1, y_1), (dx_2, dy_2)) = \int \phi(x_2) \mu^{2|1}(x_1, dx_2) \\ \int \psi(y_2) \gamma^{2|1}((x_1, y_1), (dx_2, dy_2)) = \int \psi(y_2) \nu^{2|1}(y_1, dy_2). \end{cases}$$

As we want to prove it for γ^1 -a.e. (x_1, y_1) , it is sufficient to prove that for all continuous function ξ ,

$$\begin{cases} \iint \xi(x_1, y_1) \phi(x_2) \gamma^{2|1}((x_1, y_1), (dx_2, dy_2)) d\gamma^1(x_1, y_1) = \iint \xi(x_1, y_1) \phi(x_2) \mu^{2|1}(x_1, dx_2) d\gamma^1(x_1, y_1) \\ \iint \xi(x_1, y_1) \psi(y_2) \gamma^{2|1}((x_1, y_1), (dx_2, dy_2)) d\gamma^1(x_1, y_1) = \iint \xi(x_1, y_1) \psi(y_2) \nu^{2|1}(y_1, dy_2) d\gamma^1(x_1, y_1). \end{cases}$$

First, we can use the projections $\pi_x(x, y) = x$ and $\pi_y(x, y) = y$. Moreover, we know that $\gamma^1 = (Id \times T_K^1)_\# \mu^1$. The alternate Knothe-Rosenblatt rearrangement is, as the usual one, bijective (because μ and ν are absolutely continuous), and thus, as we suppose that ν satisfies the same hypothesis than μ , we also have $\gamma^1 = ((T_K^1)^{-1}, Id)_\# \nu^1$. Let's note $\tilde{T}_K^1 = (T_K^1)^{-1}$. Then, the equalities that we want to show are

$$\begin{cases} \iint \xi(x_1, T_K^1(x_1)) \phi(x_2) \gamma_x^{2|1}((x_1, T_K^1(x_1)), dx_2) d\mu^1(x_1) = \iint \xi(x_1, T_K^1(x_1)) \phi(x_2) \mu^{2|1}(x_1, dx_2) d\mu^1(x_1) \\ \iint \xi(\tilde{T}_K^1(y_1), y_1) \psi(y_2) \gamma_y^{2|1}((\tilde{T}_K^1(y_1), y_1), dy_2) d\nu^1(y_1) = \iint \xi(\tilde{T}_K^1(y_1), y_1) \psi(y_2) \nu^{2|1}(y_1, dy_2) d\nu^1(y_1). \end{cases}$$

And we have indeed

$$\begin{aligned} \iint \xi(x_1, T_K^1(x_1)) \phi(x_2) \gamma_x^{2|1}((x_1, T_K^1(x_1)), dx_2) d\mu^1(x_1) &= \iint \xi(x_1, T_K^1(x_1)) \phi(x_2) d\gamma^{1,2}((x_1, x_2), (y_1, y_2)) \\ &= \iint \xi(x_1, T_K^1(x_1)) \phi(x_2) d\gamma_x^{1,2}(x_1, x_2) \\ &= \iint \xi(x_1, T_K^1(x_1)) \phi(x_2) \mu^{2|1}(x_1, dx_2) d\mu^1(x_1). \end{aligned}$$

We can do the same for the ν part by symmetry.

Part 3:

Now, we can proceed the same way by induction. Let $\ell \in \{2, \dots, d\}$ and suppose that the result is true in dimension $\ell - 1$ (i.e. $\gamma^{1:\ell-1} = \pi_{\#}^{1:\ell-1} \gamma = \gamma_K^{1:\ell-1}$).

For this part of the proof, we rely on (Santambrogio, 2015)[Theorem 2.23]. We can build a measure $\gamma_K^t \in \mathcal{P}(\mathbb{R}^d \times \mathbb{R}^d)$ such that

$$\begin{cases} \pi_{\#}^x \gamma_K^t = \mu \\ \pi_{\#}^y \gamma_K^t = \nu \\ \pi_{\#}^{1:\ell-1} \gamma_K^t = \eta_{t,\ell} \end{cases} \quad (16)$$

where $\eta_{t,\ell}$ is the optimal transport plan between $\mu^\ell = \pi_{\#}^{1:\ell-1} \mu$ and $\nu^\ell = \pi_{\#}^{1:\ell-1} \nu$ for the objective

$$\iint \sum_{k=1}^{\ell-1} \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma(x, y) d\gamma(x', y').$$

By induction hypothesis, we have $\eta_{t,\ell} \xrightarrow[t \rightarrow 0]{\mathcal{D}} \pi_{\#}^{1:\ell-1} \gamma_K$. To build such a measure, we can first disintegrate μ and ν

$$\begin{cases} \mu = \mu^{1:\ell-1} \otimes \mu^{\ell:d|1:\ell-1} \\ \nu = \nu^{1:\ell-1} \otimes \nu^{\ell:d|1:\ell-1}, \end{cases}$$

then we pick the Knothe transport $\gamma_K^{\ell:d|1:\ell-1}$ between $\mu^{\ell:d|1:\ell-1}$ and $\nu^{\ell:d|1:\ell-1}$. Thus, by taking $\gamma_K^T = \eta_{t,\ell} \otimes \gamma_K^{\ell:d|1:\ell-1}$, γ_K^T satisfies well the conditions (16).

Hence, we have,

$$\begin{aligned} & \iint \sum_{k=1}^{\ell-1} \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_K^t(x, y) d\gamma_K^t(x', y') \\ &= \iint \sum_{k=1}^{\ell-1} \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\eta_{t,\ell}(x_{1:\ell-1}, y_{1:\ell-1}) d\eta_{t,\ell}(x'_{1:\ell-1}, y'_{1:\ell-1}) \\ &\leq \iint \sum_{k=1}^{\ell-1} \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_t(x, y) d\gamma_t(x', y'), \end{aligned}$$

and therefore

$$\begin{aligned} & \iint \sum_{k=1}^{\ell-1} \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_K^t(x, y) d\gamma_K^t(x', y') + \iint \sum_{k=\ell}^d \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_t(x, y) d\gamma_t(x', y') \\ &\leq \mathcal{HW}_t^2(\mu, \nu) \\ &\leq \iint \sum_{k=1}^{\ell-1} \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_K^t(x, y) d\gamma_K^t(x', y') + \iint \sum_{k=\ell}^d \left(\prod_{i=1}^{k-1} \lambda_t^{(i)} \right) (x_k x'_k - y_k y'_k)^2 d\gamma_K^t(x, y) d\gamma_K^t(x', y'). \end{aligned}$$

As before, by subtracting the first term, and dividing by $\prod_{i=1}^{\ell-1} \lambda_t^{(i)}$, we get

$$\iint (x_{\ell} x'_{\ell} - y_{\ell} y'_{\ell})^2 d\gamma_t(x, y) d\gamma_t(x', y') \leq \iint (x_{\ell} x'_{\ell} - y_{\ell} y'_{\ell})^2 d\gamma_K^t(x, y) d\gamma_K^t(x', y').$$

For the right hand side, using that $\gamma_K^t = \eta_{t,\ell} \otimes \gamma_K^{\ell:d|1:\ell-1}$, we have

$$\begin{aligned} & \iint (x_{\ell} x'_{\ell} - y_{\ell} y'_{\ell})^2 d\gamma_K^t(x, y) d\gamma_K^t(x', y') \\ &= \iiint \int (x_{\ell} x'_{\ell} - y_{\ell} y'_{\ell})^2 \gamma_K^{\ell:d|1:\ell-1}((x_{1:\ell-1}, y_{1:\ell-1}), (dx_{\ell:d}, dy_{\ell:d})) \gamma_K^{\ell:d|1:\ell-1}((x'_{1:\ell-1}, y'_{1:\ell-1}), (dx'_{\ell:d}, dy'_{\ell:d})) \\ & d\eta_{t,\ell}(x_{1:\ell-1}, y_{1:\ell-1}) d\eta_{t,\ell}(x'_{1:\ell-1}, y'_{1:\ell-1}) \\ &= \iiint \int (x_{\ell} x'_{\ell} - y_{\ell} y'_{\ell})^2 \gamma_K^{\ell|1:\ell-1}((x_{1:\ell-1}, y_{1:\ell-1}), (dx_{\ell}, dy_{\ell})) \gamma_K^{\ell|1:\ell-1}((x'_{1:\ell-1}, y'_{1:\ell-1}), (dx'_{\ell}, dy'_{\ell})) \\ & d\eta_{t,\ell}(x_{1:\ell-1}, y_{1:\ell-1}) d\eta_{t,\ell}(x'_{1:\ell-1}, y'_{1:\ell-1}). \end{aligned}$$

Let's note for $\eta_{t,\ell}$ almost every $(x_{1:\ell-1}, y_{1:\ell-1}), (x'_{1:\ell-1}, y'_{1:\ell-1})$

$$GW(\mu^{\ell|1:\ell-1}, \nu^{\ell|1:\ell-1}) = \iint (x_{\ell} x'_{\ell} - y_{\ell} y'_{\ell})^2 \gamma_K^{\ell|1:\ell-1}((x_{1:\ell-1}, y_{1:\ell-1}), (dx_{\ell}, dy_{\ell})) \gamma_K^{\ell|1:\ell-1}((x'_{1:\ell-1}, y'_{1:\ell-1}), (dx'_{\ell}, dy'_{\ell})),$$

then

$$\iint (x_{\ell} x'_{\ell} - y_{\ell} y'_{\ell})^2 d\gamma_K^t(x, y) d\gamma_K^t(x', y') = \iint GW(\mu^{\ell|1:\ell-1}, \nu^{\ell|1:\ell-1}) d\eta_{t,\ell}(x_{1:\ell-1}, y_{1:\ell-1}) d\eta_{t,\ell}(x'_{1:\ell-1}, y'_{1:\ell-1}).$$

By Theorem 4, we have $\eta_{t,\ell} \otimes \eta_{t,\ell} \xrightarrow[t \rightarrow 0]{\mathcal{D}} \pi_{\#}^{1:\ell-1} \gamma_K \otimes \pi_{\#}^{1:\ell-1} \gamma_K$. So, if $\eta \mapsto \iint GW(\mu^{\ell|1:\ell-1}, \nu^{\ell|1:\ell-1}) d\eta d\eta$ is continuous over the transport plans between $\mu^{1:\ell-1}$ and $\nu^{1:\ell-1}$, we have

$$\begin{aligned} & \iint (x_{\ell} x'_{\ell} - y_{\ell} y'_{\ell})^2 d\gamma_K^t(x, y) d\gamma_K^t(x', y') \\ & \xrightarrow[t \rightarrow 0]{} \iint GW(\mu^{\ell|1:\ell-1}, \nu^{\ell|1:\ell-1}) \pi_{\#}^{1:\ell-1} \gamma_K(dx_{1:\ell-1}, dy_{1:\ell-1}) \pi_{\#}^{1:\ell-1} \gamma_K(dx'_{1:\ell-1}, dy'_{1:\ell-1}) \end{aligned}$$

and

$$\begin{aligned} & \iint GW(\mu^{\ell|1:\ell-1}, \nu^{\ell|1:\ell-1}) \pi_{\#}^{1:\ell-1} \gamma_K(dx_{1:\ell-1}, dy_{1:\ell-1}) \pi_{\#}^{1:\ell-1} \gamma_K(dx'_{1:\ell-1}, dy'_{1:\ell-1}) \\ &= \iint (x_{\ell} x'_{\ell} - y_{\ell} y'_{\ell})^2 d\gamma_K(x, y) d\gamma_K(x', y') \end{aligned}$$

by replacing the true expression of GW and using the disintegration $\gamma_K = (\pi_K^{1:\ell-1})_{\#} \gamma_K \otimes \gamma_K^{\ell|1:\ell-1}$.

For the continuity, we can apply (Santambrogio, 2015)[Lemma 1.8] (as in the (Santambrogio, 2015)[Corollary 2.24]) with $X = Y = \mathbb{R}^{\ell-1} \times \mathbb{R}^{\ell-1}$, $\tilde{X} = \tilde{Y} = \mathcal{P}(\Omega)$ with $\Omega \subset \mathbb{R}^{d-\ell+1} \times \mathbb{R}^{d-\ell+1}$ and $c(a, b) = GW(a, b)$ which can be bounded on compact supports by $\max |c|$. Moreover, we use Theorem 4 and the fact that $\eta_t \otimes \eta_t \xrightarrow[t \rightarrow 0]{\mathcal{D}} \gamma_K^{1:\ell-1} \otimes \gamma_K^{1:\ell-1}$.

By taking the limit $t \rightarrow 0$, we now get

$$\iint (x_{\ell} x'_{\ell} - y_{\ell} y'_{\ell})^2 d\gamma(x, y) d\gamma(x', y') \leq \iint (x_{\ell} x'_{\ell} - y_{\ell} y'_{\ell})^2 d\gamma_K(x, y) d\gamma_K(x', y').$$

We can now disintegrate with respect to $\gamma^{1:\ell-1}$ as before. We just need to prove that the marginals coincide which is done by taking for test functions

$$\begin{cases} \xi(x_1, \dots, x_{\ell-1}, y_1, \dots, y_{\ell-1}) \phi(x_{\ell}) \\ \xi(x_1, \dots, x_{\ell-1}, y_1, \dots, y_{\ell-1}) \psi(y_{\ell}) \end{cases}$$

and using the fact that the measures are concentrated on $y_k = T_K(x_k)$.

Part 4:

Therefore, we have well $\gamma_t \xrightarrow[t \rightarrow 0]{\mathcal{D}} \gamma_K$. Finally, for the L^2 convergence, we have

$$\int \|T_t(x) - T_K(x)\|_2^2 \mu(dx) = \int \|y - T_K(x)\|_2^2 d\gamma_t(x, y) \rightarrow \int \|y - T_K(x)\|_2^2 d\gamma_K(x, y) = 0$$

as $\gamma_t = (Id \times T_t)_{\#} \mu$ and $\gamma_K = (Id \times T_K)_{\#} \mu$. Hence, $T_t \xrightarrow[t \rightarrow 0]{L^2} T_K$.

□

C Solving \mathcal{HW} in the discrete setting

In this part, we derive formulas to solve numerically \mathcal{HW} .

Let $x_1, \dots, x_n \in \mathbb{R}^d$, $y_1, \dots, y_m \in \mathbb{R}^d$, $\alpha \in \Sigma_n$, $\beta \in \Sigma_m$, $p = \sum_{i=1}^n \alpha_i \delta_{x_i}$ and $q = \sum_{j=1}^m \beta_j \delta_{y_j}$ two discrete measures in \mathbb{R}^d . The Hadamard Wasserstein problem (7) becomes in the discrete setting

$$\begin{aligned} \mathcal{HW}^2(p, q) &= \inf_{\gamma \in \Pi(p, q)} \sum_{i,j} \sum_{k,\ell} \|x_i \odot x_k - y_j \odot y_{\ell}\|_2^2 \gamma_{i,j} \gamma_{k,\ell} \\ &= \inf_{\gamma \in \Pi(p, q)} \mathcal{E}(\gamma) \end{aligned}$$

with $\mathcal{E}(\gamma) = \sum_{i,j} \sum_{k,\ell} \|x_i \odot x_k - y_j \odot y_{\ell}\|_2^2 \gamma_{i,j} \gamma_{k,\ell}$. As denoted in Peyré et al. (2016), if we note $\mathcal{L}_{i,j,k,\ell} = \|x_i \odot x_k - y_j \odot y_{\ell}\|_2^2$, then we have

$$\mathcal{E}(\gamma) = \langle \mathcal{L} \otimes \gamma, \gamma \rangle,$$

where \otimes is defined as

$$\mathcal{L} \otimes \gamma = \left(\sum_{k,\ell} \mathcal{L}_{i,j,k,\ell} \gamma_{k,\ell} \right)_{i,j} \in \mathbb{R}^{n \times m}.$$

Proposition 7. Let $\gamma \in \Pi(p, q) = \{M \in (\mathbb{R}_+)^{n \times m}, M\mathbb{1}_m = p, M^T\mathbb{1}_n = q\}$, where $\mathbb{1}_n = (1, \dots, 1)^T \in \mathbb{R}^n$. Let's note $X = (x_i \odot x_k)_{i,k} \in \mathbb{R}^{n \times n \times d}$, $Y = (y_j \odot y_\ell)_{j,\ell} \in \mathbb{R}^{m \times m \times d}$, $X^{(2)} = (\|X_{i,k}\|_2^2)_{i,k} \in \mathbb{R}^{n \times n}$, $Y^{(2)} = (\|Y_{j,\ell}\|_2^2)_{j,\ell} \in \mathbb{R}^{m \times m}$ and, $\forall t \in \{1, \dots, d\}$, $X_t = (X_{i,k,t})_{i,k} \in \mathbb{R}^{n \times n}$ and $Y_t = (Y_{j,\ell,t})_{j,\ell} \in \mathbb{R}^{m \times m}$. Then,

$$\mathcal{L} \otimes \gamma = X^{(2)}p\mathbb{1}_m^T + \mathbb{1}_nq^T(Y^{(2)})^T - 2 \sum_{t=1}^d X_t \gamma Y_t^T.$$

Proof of Proposition 7. First, we can start by writing

$$\begin{aligned} \mathcal{L}_{i,j,k,\ell} &= \|x_i \odot x_k - y_j \odot y_\ell\|_2^2 \\ &= \|X_{i,k} - Y_{j,\ell}\|_2^2 \\ &= \|X_{i,k}\|_2^2 + \|Y_{j,\ell}\|_2^2 - 2\langle X_{i,k}, Y_{j,\ell} \rangle \\ &= [X^{(2)}]_{i,k} + [Y^{(2)}]_{j,\ell} - 2\langle X_{i,k}, Y_{j,\ell} \rangle. \end{aligned}$$

We cannot directly apply proposition 1 from (Peyré et al., 2016) (as the third term is a scalar product), but by doing the same type of computation, we get

$$\mathcal{L} \otimes \gamma = A + B + C$$

with

$$A_{i,j} = \sum_{k,\ell} [X^{(2)}]_{i,k} \gamma_{k,\ell} = \sum_k [X^{(2)}]_{i,k} \sum_\ell \gamma_{k,\ell} = \sum_k [X^{(2)}]_{i,k} [\gamma \mathbb{1}_m]_{k,1} = [X^{(2)} \gamma \mathbb{1}_m]_{i,1} = [X^{(2)} p]_{i,1}$$

$$B_{i,j} = \sum_{k,\ell} [Y^{(2)}]_{j,\ell} \gamma_{k,\ell} = \sum_\ell [Y^{(2)}]_{j,\ell} \sum_k \gamma_{k,\ell} = \sum_\ell [Y^{(2)}]_{j,\ell} [\gamma^T \mathbb{1}_n]_{\ell,1} = [Y^{(2)} \gamma^T \mathbb{1}_n]_{j,1} = [Y^{(2)} q]_{j,1}$$

and

$$\begin{aligned} C_{i,j} &= -2 \sum_{k,\ell} \langle X_{i,k}, Y_{j,\ell} \rangle \gamma_{k,\ell} = -2 \sum_{k,\ell} \sum_{t=1}^d X_{i,k,t} Y_{j,\ell,t} \gamma_{k,\ell} \\ &= -2 \sum_{t=1}^d \sum_k [X_t]_{i,k} \sum_\ell [Y_t]_{j,\ell} \gamma_{\ell,k}^T \\ &= -2 \sum_{t=1}^d \sum_k [X_t]_{i,k} [Y_t \gamma^T]_{j,k} \\ &= -2 \sum_{t=1}^d [X_t (Y_t \gamma^T)^T]_{i,j}. \end{aligned}$$

Finally, we have

$$\mathcal{L} \otimes \gamma = X^{(2)}p\mathbb{1}_m^T + \mathbb{1}_nq^T(Y^{(2)})^T - 2 \sum_{t=1}^d X_t \gamma Y_t^T.$$

□

Remark 1. The complexity for computing $\mathcal{L} \otimes \gamma$ is $O(d(n^2m + m^2n))$.

Remark 2. For the degenerated cost function (8), we just need to replace X and Y by $\tilde{X}_t = A_t^{\frac{1}{2}} X$ and $\tilde{Y}_t = A_t^{\frac{1}{2}} Y$ in the previous proposition.

To solve this problem numerically, we can use the conditional gradient algorithm (Algorithm 2 in (Vayer et al., 2019a)). This algorithm only requires to compute the gradient

$$\nabla \mathcal{E}(\gamma) = 2(A + B + C) = 2(\mathcal{L} \otimes \gamma)$$

at each step and a classical OT problem. This algorithm is more efficient than solving directly the quadratic problem. Moreover, while it is a non convex problem, it actually converges to a local stationary point (Lacoste-Julien, 2016).

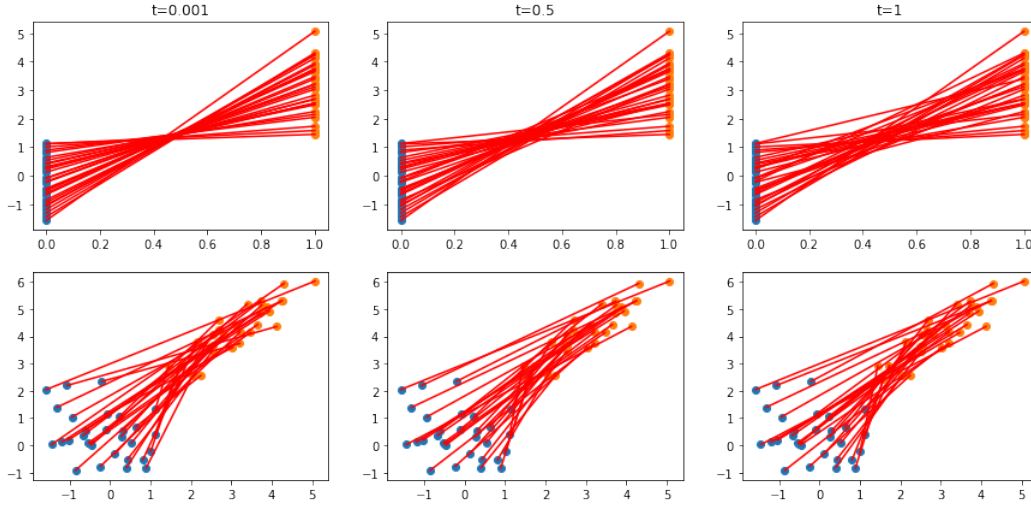


Figure 3: Degenerated Coupling

On Figure 3, we generated 30 points of two gaussian distributions, and computed the optimal coupling of \mathcal{HW}_t for several t . These points have the same uniform weight. On the first row, we projected the points on the first coordinate. Note that for discrete points, the Knothe-Rosenblatt coupling comes back to sort the point with respect to the first coordinate if there is no ambiguity (*i.e.* $x_1^{(1)} < \dots < x_n^{(1)}$) as it comes back to perform the optimal transport in one dimension (Peyré et al., 2019)[Remark 2.28]. For our cost, the optimal coupling in 1D can either be the increasing or the decreasing one. We observe well on the first row of figure (3) that the optimal coupling when t is close to 0 corresponds to the “anti-cdf”.