



HAL
open science

The DTM-signature for a geometric comparison of metric-measure spaces from samples

Claire Brécheteau

► **To cite this version:**

Claire Brécheteau. The DTM-signature for a geometric comparison of metric-measure spaces from samples. 2017. hal-01426331v2

HAL Id: hal-01426331

<https://inria.hal.science/hal-01426331v2>

Preprint submitted on 9 Feb 2017 (v2), last revised 20 Feb 2019 (v3)

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.

The DTM-signature for a geometric comparison of metric-measure spaces from samples*

Claire Bréchet

Université Paris-Sud – Inria Saclay // Université Paris-Saclay, France
claire.brecheteau@inria.fr

February 9, 2017

Abstract

In this paper, we introduce the notion of DTM-signature, a measure on \mathbb{R}_+ that can be associated to any metric-measure space. This signature is based on the distance to a measure (DTM) introduced by Chazal, Cohen-Steiner and Mérigot. It leads to a pseudo-metric between metric-measure spaces, upper-bounded by the Gromov-Wasserstein distance. Under some geometric assumptions, we derive lower bounds for this pseudo-metric.

Given two N -samples, we also build an asymptotic statistical test based on the DTM-signature, to reject the hypothesis of equality of the two underlying metric-measure spaces, up to a measure-preserving isometry. We give strong theoretical justifications for this test and propose an algorithm for its implementation.

1 Introduction

Among the variety of data available, from astrophysics to biology, including social networks and so on, many come as sets of points from a metric space. A natural question, given two sets of such data is to decide whether they are similar, that is whether they come from the same distribution, whether their shape are close, or not. This comparison may be compromised when the data are not embedded into the same space, or if the two systems of coordinates in which the data are represented are different. To overcome this issue, a natural idea is to forget about this embedding and only consider the set of points together with the distances between pairs. A natural framework to compare data is then to assume that they come from a measure on a metric space and to consider two such metric-measure spaces as being the same when they are equal up to some isomorphism, as defined below.

Definition 1 (mm-space).

A *metric-measure space* (**mm-space**) is a triple $(\mathcal{X}, \delta, \mu)$, with \mathcal{X} a set, δ a metric on \mathcal{X} and μ a probability measure on \mathcal{X} equipped with its Borel σ -algebra.

Definition 2 (Isomorphism between mm-spaces).

Two mm-spaces $(\mathcal{X}, \delta, \mu)$ and $(\mathcal{Y}, \gamma, \nu)$ are said to be **isomorphic** if there exist Borel sets $\mathcal{X}_0 \subset \mathcal{X}$ and $\mathcal{Y}_0 \subset \mathcal{Y}$ such that $\mu(\mathcal{X} \setminus \mathcal{X}_0) = 0$ and $\nu(\mathcal{Y} \setminus \mathcal{Y}_0) = 0$, and some one-to-one and onto isometry $\phi : \mathcal{X}_0 \rightarrow \mathcal{Y}_0$ preserving measures, that is satisfying $\nu(\phi(A \cap \mathcal{X}_0)) = \mu(A \cap \mathcal{X}_0)$ for any Borel set A of \mathcal{X} . Such a map ϕ is called an **isomorphism** between the mm-spaces $(\mathcal{X}, \delta, \mu)$ and $(\mathcal{Y}, \gamma, \nu)$.

*This work was partially supported by the ANR project TopData and GUDHI

In this paper, we first address the question of the comparison of general mm-spaces, up to an isomorphism. In other terms, we aim at designing a metric or at least a pseudo-metric on the quotient space of mm-spaces by the relation of isomorphism. A suitable pseudo-distance should be stable under some perturbations, under sampling, discriminative and easy to implement when dealing with discrete spaces.

A first characterisation of mm-spaces is given in [19]. In its Theorem 3 $\frac{1}{2}$.5, Gromov proves that any mm-space can be recovered, up to an isomorphism, from the knowledge, for all size N , of the distribution of the $N \times N$ -matrix of distances associated to a N -sample. More recently, in [23], Mémoli proposes metrics on the quotient space of mm-spaces by the relation of isomorphism, the Gromov–Wasserstein distances.

Definition 3 (Gromov–Wasserstein distance).

The **Gromov–Wasserstein distance** between two mm-spaces $(\mathcal{X}, \delta, \mu)$ and $(\mathcal{Y}, \gamma, \nu)$ with parameter $p \in [1, \infty)$ denoted $GW_p(\mathcal{X}, \mathcal{Y})$ is defined by the expression:

$$\inf_{\pi \in \Pi(\mu, \nu)} \frac{1}{2} \left(\int_{\mathcal{X} \times \mathcal{Y}} \int_{\mathcal{X} \times \mathcal{Y}} (\Gamma_{\mathcal{X}, \mathcal{Y}}(x, y, x', y'))^p \pi(\mathrm{d}x \times \mathrm{d}y) \pi(\mathrm{d}x' \times \mathrm{d}y') \right)^{\frac{1}{p}},$$

with $\Gamma_{\mathcal{X}, \mathcal{Y}}(x, y, x', y') = |\delta(x, x') - \gamma(y, y')|$. Here $\Pi(\mu, \nu)$ stands for the set of **transport plans** between μ and ν , that is the set of Borel probability measures π on $\mathcal{X} \times \mathcal{Y}$ satisfying $\pi(A \times \mathcal{Y}) = \mu(A)$ and $\pi(\mathcal{X} \times B) = \nu(B)$ for all Borel sets A in \mathcal{X} and B in \mathcal{Y} .

Unfortunately, even when dealing with discrete mm-spaces, the computation of these Gromov–Wasserstein distances is extremely costly. An alternative is to build a **signature** from each mm-space, that is an object invariant under isomorphism. The mm-spaces are then compared through their signatures. In [23], Mémoli gives an overview of such signatures, as for instance shape distribution, eccentricity or what he calls local distribution of distances.

In this paper, we introduce a new signature that is a probability measure on \mathbb{R} , and we propose to compare such signatures using Wasserstein distances [26].

Definition 4 (Wasserstein distance).

The **Wasserstein distance** of parameter $p \in [1, \infty)$ between two Borel probability measures μ and ν over the same metric space (\mathcal{X}, δ) is defined as:

$$W_p(\mu, \nu) = \inf_{\pi \in \Pi(\mu, \nu)} \left(\int_{\mathbb{R}^2} \delta^p(x, y) \mathrm{d}\pi(x, y) \right)^{\frac{1}{p}}.$$

For two probability measures μ and ν over \mathbb{R}_+ , the L_1 -Wasserstein distance can be rewritten as the L_1 -norm between the cumulative distribution functions of the measures, $F_\mu : t \mapsto \mu((-\infty, t])$ and F_ν , or as well, as the L_1 -norm between the quantile functions, $F_\mu^{-1} : s \mapsto \inf\{x \in \mathbb{R} \mid F_\mu(x) \geq s\}$ and F_ν^{-1} . Thus, the computation of the L_1 -Wasserstein distance between empirical measures is easy, in $O(N \log(N))$ for two empirical measures from subsets of \mathbb{R} of size N , the complexity of a sort.

Shape signatures are widely used for classification or pre-classification tasks; see for instance [25]. With a more topological point of view, persistence diagrams have been used for this purpose in [10, 12]. But, as far as we know, the construction of well-founded statistical tests from signatures to compare mm-spaces has not been considered among the literature. This is the second problem focussed in this paper.

Recall that a **statistical test** is a random variable ϕ_N taking values in $\{0, 1\}$. More precisely ϕ_N is a function of N random data from a distribution \mathcal{L}_θ depending on some unknown parameter θ in some set Θ . It is associated to two hypotheses H_0 “ $\theta \in \Theta_0$ ” and H_1 “ $\theta \in \Theta_1$ ” with Θ_0 and Θ_1 disjoint subsets of Θ . Ideally, we would like the test ϕ_N to be equal to 1 if θ is in Θ_1 and to be 0 if θ is in Θ_0 .

The quality of a statistical test is measured in terms of its **type I error**, that is the function defined for all θ_0 in Θ_0 by $\mathbb{P}_{\theta_0}(\phi_N = 1)$, the probability of pretending θ to be in Θ_1 when $\theta = \theta_0$ is actually in Θ_0 . Moreover, a test is of **level** $\alpha \in (0, 1)$ if its type I error is upper-bounded by α , that is $\mathbb{P}_{\theta_0}(\phi_N = 1) \leq \alpha$ for all θ_0 in Θ_0 . Two statistical tests with a fixed level $\alpha \in (0, 1)$ can be compared through their **type II error**, that is the function defined for all θ_1 in Θ_1 by $\mathbb{P}_{\theta_1}(\phi_N = 0)$, the probability of pretending θ to be in Θ_0 when $\theta = \theta_1$ is actually in Θ_1 . See [4] for a reference on statistical tests.

In this article, we build a **test of asymptotic level** α , that is a test ϕ_N such that for any θ_0 in Θ_0 , $\mathbb{P}_{\theta_0}(\phi_N = 1) \rightarrow \alpha$ when the size of the sample N goes to ∞ . Moreover, the set Θ we consider is the set of couples of mm-spaces $((\mathcal{X}, \delta, \mu), (\mathcal{Y}, \gamma, \nu))$. The set Θ_0 is the subset of Θ made of couples of two isomorphic spaces: $\Theta_0 = \{((\mathcal{X}, \delta, \mu), (\mathcal{Y}, \gamma, \nu)) \in \Theta \mid (\mathcal{X}, \delta, \mu) \text{ and } (\mathcal{Y}, \gamma, \nu) \text{ are isomorphic}\}$, and $\Theta_1 = \Theta \setminus \Theta_0$.

Such a test generalises two-sample tests, from the precursor Kolmogorov-Smirnov test to the more recent tests in [18] or [9]. Our test does not depend on the embedding of the data and keeps a track of the geometry in some way, a point of view that has already been taken in the context of density estimation [21]. Thus, it could be of interest for proteins, 3D-shape comparison, etc.

Concretely, in this paper, we propose a new signature based on the distance to a measure (DTM) introduced in [11], the **DTM-signature**. This signature is invariant under isomorphism and easy to compute. We prove its stability with respect to the so-called Gromov-Wasserstein and Wasserstein distances with parameter $p = 1$. It leads to a stability under sampling, at least for the Euclidean space \mathbb{R}^d . After deriving frameworks under which the knowledge of the distance to a measure determines the measure, we prove discriminative properties for the DTM-signature by deriving lower bounds for the L_1 -Wasserstein distance between two such signatures, under various assumptions. Finally, from two N -samples, we derive a statistical test, based on bootstrap methods, to reject or not the hypothesis of equality of the two underlying metric-measure spaces, up to a measure-preserving isometry. This test comes with an easy-to-implement algorithm, and a strong theoretical justification.

The DTM-signature depends on some parameter $m \in (0, 1)$. It thus offers a variety of new figures, as well as new lower-bounds for the Gromov-Wasserstein distance. As for the statistical test, it presents the advantage of not depending on the embedding of the data, only the knowledge of the distances between points is required. In this sense, it is new. The justification of the validity of the test with the use of the Wasserstein distance is quite new as well, and still poorly used; see [14] for another use.

The paper is organized as follows. Section 2 is devoted to the distance to a measure. An accent is put on its discriminative properties. The DTM-signature is then introduced in Section 3. The question of discrimination of two mm-spaces is also discussed. For this purpose, we derive lower bounds for our pseudo-distance, the L_1 -Wasserstein distance between the two DTM-signatures. Finally, in Section 4 we introduce the test of isomorphism, propose an algorithm for its implementation and then give some theoretical results to ensure the validity of the procedure. Numerical illustrations are given in Section 5.

2 The distance to a measure to discriminate between measures

Let (\mathcal{X}, δ) be a metric space, equipped with a Borel probability measure μ . Given m in $[0, 1]$, the **pseudo-distance function** is defined at any point x of \mathcal{X} , by:

$$\delta_{\mu, m}(x) = \inf\{r > 0 \mid \mu(\bar{B}(x, r)) > m\}.$$

The function **distance to the measure** μ with mass parameter m and denoted $d_{\mu,m}$ is then defined for all x in \mathcal{X} by:

$$d_{\mu,m}(x) = \frac{1}{m} \int_{l=0}^m \delta_{\mu,l}(x) dl.$$

The distance to a measure is a generalisation of the function distance to a compact set; see [11]. This function is continuous with respect to the mass parameter m , and Lipschitz with respect to μ .

Proposition 5 (Stability, in [11] for \mathbb{R}^d , in [7] for metric spaces).

For two m -spaces $(\mathcal{X}, \delta, \mu)$ and $(\mathcal{Y}, \delta, \nu)$ embedded into the same metric space, we have that

$$\|d_{\mu,m} - d_{\nu,m}\|_{\infty, \mathcal{X} \cup \mathcal{Y}} \leq \frac{1}{m} W_1(\mu, \nu).$$

Moreover, for some empirical measure $\hat{\mu}_N = \frac{1}{N} \sum_{i=1}^N \delta_{X_i}$ on a metric space (\mathcal{X}, δ) , the distance to the measure $\hat{\mu}_N$ with mass parameter $\frac{k}{N}$ for some k in $\llbracket 0, N \rrbracket$ at a point x of \mathcal{X} satisfies:

$$d_{\hat{\mu}_N, \frac{k}{N}}(x) = \frac{1}{k} \sum_{i=1}^k \delta(X^{(i)}, x),$$

where $X^{(1)}, X^{(2)}, \dots, X^{(k)}$ are k nearest neighbours of x among the N points X_1, X_2, \dots, X_N .

The distance to the measure $\hat{\mu}_N$ is thus equal to the mean of the distances to k -nearest neighbours. In particular, in this case, the computation of the DTM boils down to the computation of the first k -nearest neighbours.

The question of determining if the knowledge of the distance to a measure leads to the knowledge of the measure itself is a natural question. Some work has been done in this direction for discrete measures; see [7]. In the following, we propose results in different settings.

Proposition 6.

Let (\mathcal{X}, δ) be a metric space, and $\mathcal{M}_1(\mathcal{X})$ be the set of Borel probability measures over (\mathcal{X}, δ) . We define the maps ϕ and ψ for all μ in $\mathcal{M}_1(\mathcal{X})$ by:

$$\phi(\mu) = (d_{\mu,m}(x))_{m \in [0,1], x \in \mathcal{X}}$$

and

$$\psi(\mu) = (\mu(\overline{\mathbb{B}}(x, r)))_{r \in \mathbb{R}_+, x \in \mathcal{X}}.$$

Then, the map ϕ is injective if and only if the map ψ is injective.

Proof

From the definition of $d_{\mu,m}(x)$, we have:

$$\mu(\overline{\mathbb{B}}(x, r)) = \inf\{m \geq 0 \mid \delta_{\mu,m}(x) > r\}.$$

Moreover, since $m \rightarrow \delta_{\mu,m}(x)$ is right-continuous, after the differentiation the distance-to-a-measure function with respect to m , we have:

$$m \frac{\partial}{\partial m} d_{\mu,m}(x) + d_{\mu,m}(x) = \delta_{\mu,m}(x).$$

■

It means that in spaces on which measures are determined by their values on balls, the measures are determined by the knowledge of the distance-to-a-measure functions for all parameters m in $[0, 1]$, on all x in \mathcal{X} . Remark that the Euclidean space \mathbb{R}^d satisfies such a condition, but this is not the case of every metric space, as explained in [8].

Under the following specific framework, we will establish a stronger identifiability result.

For O a non-empty bounded open subset of \mathbb{R}^d , we define the **uniform measure** μ_O for all Borel set A of \mathbb{R}^d , by:

$$\mu_O(A) = \frac{\text{Leb}_d(O \cap A)}{\text{Leb}_d(O)},$$

with Leb_d the Lebesgue measure on \mathbb{R}^d .

We also define the **medial axis** of O , $\mathcal{M}(O)$ as the set of points in O having at least two projections onto ∂O . That is,

$$\mathcal{M}(O) = \{y \in O \mid \exists x', x'' \in \partial O, x' \neq x'', \|y - x'\|_2 = \|y - x''\|_2 = d(y, \partial O)\},$$

with $d(y, \partial O) = \inf\{\|x - y\|_2 \mid x \in \partial O\}$.

Its **reach**, $\text{Reach}(O)$, is the distance between its boundary ∂O and its medial axis $\mathcal{M}(O)$. That is,

$$\text{Reach}(O) = \inf\{\|x - y\|_2 \mid x \in \partial O, y \in \mathcal{M}(O)\}.$$

If K is a compact subset of \mathbb{R}^d , it is standard to define its reach as $\text{Reach}(K^c)$, the reach of its complement in \mathbb{R}^d . See [16] to get more familiar with these notions.

Proposition 7.

Let O and O' be two non-empty bounded open subsets of \mathbb{R}^d with positive reach, such that $O = (\overline{O})^\circ$ and $O' = (\overline{O'})^\circ$. Let m be some positive constant satisfying

$$m \leq \min(\text{Reach}(O)^d, \text{Reach}(O')^d) \frac{\omega_d}{\text{Leb}_d(O)},$$

with $\omega_d = \text{Leb}_d(B(0, 1))$, the Lebesgue volume of the unit d -dimensional ball. If for all x in \mathbb{R}^d

$$d_{\mu_O, m}(x) = d_{\mu_{O'}, m}(x),$$

then $\mu_O = \mu_{O'}$.

Proof

This is a straightforward consequence of Proposition 26, in the Appendix. The proof relies on the fact that the set of points in \mathbb{R}^d minimizing the distance to the measure μ_O is equal to $\{x \in O \mid d_{\partial O}(x) \geq \epsilon(m, O)\}$ with $\epsilon(m, O) = \left(\frac{m \text{Leb}_d(O)}{\omega_d}\right)^{\frac{1}{d}}$, providing that the set is non-empty. Then, if $\text{Reach}(O)$ is not smaller than $\epsilon(m, O)$, O equals to the set of points at distance smaller than $\epsilon(m, O)$ from $\{x \in O \mid d_{\partial O}(x) \geq \epsilon(m, O)\}$. Thus, the measure μ_O can be recovered. We use the notion of **skeleton** in [20] for some details in the proof. ■

It means that for m small enough, the knowledge of the distance to a measure at any point x in \mathbb{R}^d for two measures μ_O and $\mu_{O'}$ is discriminative.

3 The DTM-signature to discriminate between metric-measure spaces

From the distance-to-a-measure function, we derive a new signature.

Definition 8 (DTM-signature).

The **DTM-signature** associated to some mm-space $(\mathcal{X}, \delta, \mu)$, denoted $d_{\mu,m}(\mu)$, is the distribution of the real-valued random variable $d_{\mu,m}(X)$ where X is some random variable of law μ .

The DTM-signature turns out to be stable in the following sense.

Proposition 9.

We have that:

$$W_1(d_{\mu,m}(\mu), d_{\nu,m}(\nu)) \leq \frac{1}{m} GW_1(\mathcal{X}, \mathcal{Y}).$$

Proof

Proof in the Appendix, in Section B. The proof is relatively similar to the ones given by Mémoli in [23] for other signatures. ■

It follows directly that two isomorphic mm-spaces have the same DTM-signature. Whenever the two mm-spaces are embedded into the same metric space, we also get stability with respect to the L_1 -Wasserstein distance.

Proposition 10.

If $(\mathcal{X}, \delta, \mu)$ and $(\mathcal{Y}, \delta, \nu)$ are two metric spaces embedded into some metric space (\mathcal{Z}, δ) , then we can upper bound $W_1(d_{\mu,m}(\mu), d_{\nu,m}(\nu))$ by

$$W_1(\mu, \nu) + \min \{ \|d_{\mu,m} - d_{\nu,m}\|_{\infty, \text{Supp}(\mu)}, \|d_{\mu,m} - d_{\nu,m}\|_{\infty, \text{Supp}(\nu)} \},$$

and more generally by

$$\left(1 + \frac{1}{m}\right) W_1(\mu, \nu).$$

Proof

First remark that:

$$\begin{aligned} W_1(d_{\mu,m}(\mu), d_{\nu,m}(\mu)) &\leq \int_{\mathcal{X}} |d_{\mu,m}(x) - d_{\nu,m}(x)| d\mu(x) \\ &\leq \|d_{\mu,m} - d_{\nu,m}\|_{\infty, \text{Supp}(\mu)}. \end{aligned}$$

Then, for all π in $\Pi(\mu, \nu)$:

$$W_1(d_{\nu,m}(\mu), d_{\nu,m}(\nu)) \leq \int_{\mathcal{X} \times \mathcal{Y}} |d_{\nu,m}(x) - d_{\nu,m}(y)| d\pi(x, y).$$

Thus, since $d_{\nu,m}$ is 1-Lipschitz:

$$W_1(d_{\nu,m}(\mu), d_{\nu,m}(\nu)) \leq W_1(\mu, \nu).$$

We use Proposition 5 to conclude. ■

The DTM-signature is stable but unfortunately does not always discriminates between mm-spaces. Indeed, in the following counter-example from [23] (example 5.6), there are two non-isomorphic mm-spaces sharing the same signatures for all values of m .

Example 11. We consider two graphs made of 9 vertices each, clustered in three groups of 3 vertices, such that each vertex is at distance 1 exactly to each vertex of its group and at distance 2 to any other vertex. We assign a mass to each vertex, the distribution is the following, for the first graph:

$$\mu = \left\{ \left(\frac{23}{140}, \frac{1}{105}, \frac{67}{420} \right), \left(\frac{3}{28}, \frac{1}{28}, \frac{4}{21} \right), \left(\frac{2}{15}, \frac{1}{15}, \frac{2}{15} \right) \right\},$$

and for the second graph:

$$\nu = \left\{ \left(\frac{3}{28}, \frac{1}{15}, \frac{67}{420} \right), \left(\frac{2}{15}, \frac{4}{21}, \frac{1}{105} \right), \left(\frac{23}{140}, \frac{2}{15}, \frac{1}{28} \right) \right\}.$$

The mm-spaces ensuing are not isomorphic since any one-to-one and onto measure-preserving map would send at least one couple of vertices at distance 1 to each other, to a couple of vertices at distance 2 to each other, thus it would not be an isometry.

Moreover, remark that the DTM-signatures associated to the graphs are equal since the total mass of each cluster is exactly equal to $\frac{1}{3}$.

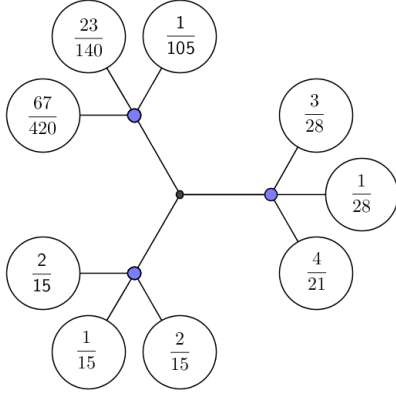


Figure 1: μ

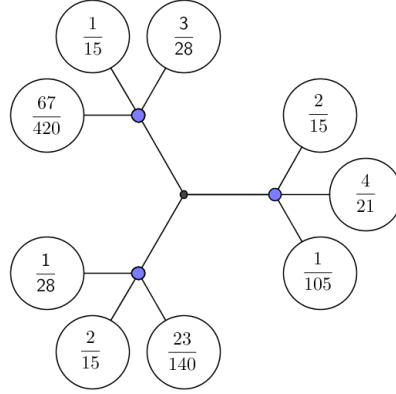


Figure 2: ν

Nevertheless, the signature can be discriminative in some cases. In the following, we give lower bounds for the L_1 -Wasserstein distance between two signatures under three different alternatives.

3.1 When the distances are multiplied by some positive real number λ

Let λ be some positive real number. The DTM-signature discriminates between two mm-spaces isomorphic up to a dilatation of parameter λ , for $\lambda \neq 1$.

Proposition 12.

Let $(\mathcal{X}, \delta, \mu)$ and $(\mathcal{Y}, \gamma, \nu) = (\mathcal{X}, \lambda\delta, \mu)$ be two mm-spaces. We have

$$W_1(d_{\mu,m}(\mu), d_{\nu,m}(\nu)) = |1 - \lambda| \mathbb{E}_\mu[d_{\mu,m}(X)],$$

for X a random variable of law μ .

Proof

First remark that $F_{d_{\nu,m}(\nu)}^{-1} = \lambda F_{d_{\mu,m}(\mu)}^{-1}$. Then,

$$\begin{aligned} W_1(d_{\mu,m}(\mu), d_{\nu,m}(\nu)) &= \int_0^1 \left| F_{d_{\mu,m}(\mu)}^{-1}(s) - F_{d_{\nu,m}(\nu)}^{-1}(s) \right| ds \\ &= |1 - \lambda| \int_0^1 \left| F_{d_{\mu,m}(\mu)}^{-1}(s) \right| ds \\ &= |1 - \lambda| \mathbb{E}_\mu [d_{\mu,m}(X)]. \end{aligned}$$

■

3.2 The case of uniform measures on non-empty bounded open subsets of \mathbb{R}^d

The DTM-signature discriminates between two uniform measures over two non-empty bounded open subsets of \mathbb{R}^d with different Lebesgue volume.

Proposition 13.

Let $(O, \|\cdot\|_2, \mu_O)$ and $(O', \|\cdot\|_2, \mu_{O'})$ be two mm-spaces, for O and O' two non-empty bounded open subsets of \mathbb{R}^d satisfying $O = (\overline{O})^\circ$ and $O' = (\overline{O'})^\circ$, and $\|\cdot\|_2$ the euclidean norm. A lower bound for $W_1(d_{\mu_O,m}(\mu_O), d_{\mu_{O'},m}(\mu_{O'}))$ is given by:

$$\min(\mu_O(O_{\epsilon(m,O)}), \mu_{O'}(O'_{\epsilon(m,O')})) \frac{d}{d+1} \left(\frac{m}{\omega_d} \right)^{\frac{1}{d}} \left| \text{Leb}_d(O)^{\frac{1}{d}} - \text{Leb}_d(O')^{\frac{1}{d}} \right|.$$

Here, $O_\epsilon = \{x \in O \mid d(x, \partial O) \geq \epsilon\}$, and $\epsilon(m, O) = \left(\frac{m \text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}}$ is the radius of any ball of μ_O mass m , included in O .

Proof

If the set $O_{\epsilon(m,O)}$ is non-empty, then the minimal value of the distance to a measure is given by:

$$\min_{x \in \mathbb{R}^d} (d_{\mu_O,m}(x)) = d_{\min} := \frac{d}{d+1} \left(\frac{m \text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}}.$$

Moreover, the points at minimal distance are exactly the points of $O_{\epsilon(m,O)}$. This is Proposition 25 in the Appendix. So, $F_{d_{\mu_O,m}(\mu_O)}(d_{\min}) = \mu_O(O_{\epsilon(m,O)})$. To conclude, we use the definition of the L_1 -Wasserstein distance as the L_1 -norm between the cumulative distribution functions. ■

3.3 The case of two measures on the same open subset of \mathbb{R}^d with one measure uniform

Let $(O, \|\cdot\|_2, \mu_O)$ and $(O, \|\cdot\|_2, \nu)$ be two mm-spaces with O a non-empty bounded open subset of \mathbb{R}^d and ν a measure absolutely continuous with respect to μ_O . Thanks to the Radon-Nikodym theorem, there is some μ_O -measurable function f on O such that for all Borel set A in O :

$$\nu(A) = \int_A f(\omega) d\mu_O(\omega).$$

We can consider the λ -super-level sets of the function f denoted by $\{f \geq \lambda\}$. As for the previous part, we will denote by $\{f \geq \lambda\}_\epsilon$ the set of points belonging to $\{f \geq \lambda\}$ whose distance to $\partial\{f \geq \lambda\}$ is at least ϵ .

Then we get the following lower bound for the L_1 -Wasserstein distance between the two signatures:

Proposition 14.

Under these hypotheses, a lower bound for $W_1(d_{\mu_O, m}(\mu_O), d_{\nu, m}(\nu))$ is given by:

$$\frac{1}{1+d} \frac{1}{\text{Leb}_d(O)} \left(\frac{m \text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}} \int_{\lambda=1}^{\infty} \frac{1}{\lambda^{\frac{1}{d}}} \max_{\lambda' \geq \lambda} \text{Leb}_d \left(\{f \geq \lambda'\} \left(\frac{m}{\omega_d} \frac{\text{Leb}_d(O)}{\lambda'} \right)^{\frac{1}{d}} \right) d\lambda.$$

Proof

Proof in the Appendix, in Section A.2. ■

When the density f is Hölder

We assume that f is Hölder on O , with positive parameters $\chi \in (0, 1]$ and $L > 0$, that is:

$$\forall x, y \in O, |f(x) - f(y)| \leq L \|x - y\|_2^\chi.$$

We also assume that $\text{Reach}(O) > 0$. Then for m small enough, the DTM-signature is discriminative.

Proposition 15.

Under the previous assumptions, if one of the following conditions is satisfied, then the quantity $W_1(d_{\mu_O, m}(\mu_O), d_{\nu, m}(\nu))$ is positive:

$$m < \frac{\omega_d}{\text{Leb}_d(O)} \min \left\{ \text{Reach}(O)^d, \left(\frac{\|f\|_{\infty, O} - 1}{2L} \right)^{\frac{d}{\chi}} \right\};$$

$$m \in \left[\frac{\omega_d}{\text{Leb}_d(O)} (\text{Reach}(O))^d, (\|f\|_{\infty, O} - 2L (\text{Reach}(O))^\chi) (\text{Reach}(O))^d \frac{\omega_d}{\text{Leb}_d(O)} \right);$$

$$m \in \left[\frac{\omega_d}{\text{Leb}_d(O)} \left(\frac{d}{\chi} \right)^{\frac{d}{\chi}} (2L)^{-\frac{d}{\chi}}, \min \left\{ m_0, \frac{\omega_d}{\text{Leb}_d(O)} (\text{Reach}(O))^{d+\chi} \frac{\chi}{d} 2L \right\} \right),$$

with $m_0 = \|f\|_{\infty, O}^{\frac{d}{\chi}+1} \frac{\omega_d}{\text{Leb}_d(O)} \left(\frac{d}{\chi} \right)^{\frac{d}{\chi}} (2L)^{-\frac{d}{\chi}} \left(\frac{\chi}{d+\chi} \right)^{\frac{\chi}{d+\chi}}$.

Moreover, under any of these conditions, we get the lower bound for the quantity $W_1(d_{\mu_O, m}(\mu_O), d_{\nu, m}(\nu))$:

$$\frac{1}{1+d} \left(\frac{m \text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}} \int_{\lambda=1}^{\infty} \frac{1}{\lambda^{1+\frac{1}{d}}} \sup_{\lambda' \geq \lambda} \nu(\{f \geq \lambda' + L\epsilon(\lambda')^\chi\} \cap O_{\epsilon(\lambda')}) d\lambda,$$

with $\epsilon(\lambda') = \lambda'^{-\frac{1}{d}} \left(\frac{m \text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}}$.

Proof

Proof in the Appendix, in Section A.2. ■

The previous examples provide several relevant cases where the DTM-signature turns out to be discriminative. It is thus appealing to use it as a tool to compare mm-spaces up to isomorphism.

4 An algorithm to compare metric-measure spaces from samples

In this section, $(\mathcal{X}, \delta, \mu)$ and $(\mathcal{Y}, \gamma, \nu)$ are two mm-spaces. We build a test of the null hypothesis

$$H_0 \text{ "The mm-spaces } (\mathcal{X}, \delta, \mu) \text{ and } (\mathcal{Y}, \gamma, \nu) \text{ are isomorphic",}$$

against its alternative:

$$H_1 \text{ "The mm-spaces } (\mathcal{X}, \delta, \mu) \text{ and } (\mathcal{Y}, \gamma, \nu) \text{ are not isomorphic"}$$

4.1 The algorithm

The test we propose is based on the fact that the DTM-signatures associated to two isomorphic mm-spaces are equal. If so, it leads to a pseudo-distance $W_1(d_{\mu,m}(\mu), d_{\nu,m}(\nu))$ equal to zero.

Let consider, in this part, a N -sample P from the measure μ , and a N -sample Q from the measure ν . A natural idea for a test is to approximate the pseudo-distance by the statistic $W_1(d_{\mathbb{1}_P,m}(\mathbb{1}_P), d_{\mathbb{1}_Q,m}(\mathbb{1}_Q))$, where $\mathbb{1}_P$ is the uniform probability measure on the set P , and to reject the hypothesis H_0 if this statistic is larger than some critical value. The choice of the critical value should rely on some parameter $\alpha \in (0, 1)$ and lead to a level α for the test. It strongly depends on the measures μ and ν that are unknown. Nonetheless, there exist classical ways of approximating a critical value, one is to mimic the distribution of the statistic by replacing the distribution μ with the distribution $\mathbb{1}_P$ and ν with $\mathbb{1}_Q$. Unfortunately, this standard method known as **bootstrap** fails theoretically and experimentally for our framework.

Thus, we propose another kind of bootstrap. For this purpose, we need to take P' a subset of P and Q' a subset of Q . The statistic we focus on is $W_1(d_{\mathbb{1}_{P'},m}(\mathbb{1}_{P'}), d_{\mathbb{1}_{Q'},m}(\mathbb{1}_{Q'}))$. It turns out that in this case, the critical value associated to this statistic can be well approximated from the samples P and Q , for a suitable size of P' and Q' with respect to N .

This approach leads to the following algorithm.

Algorithm 1: Test Procedure

```

Input :  $P$  and  $Q$   $N$ -samples from  $\mu$  (respectively  $\nu$ ),  $N$ ,  $n$ ,  $m$ ,  $\alpha$ ,  $N_{MC}$  even ;
# Compute  $T$  the test statistic
Take  $P'$  a random subset of  $P$  of size  $n$  ;
Take  $Q'$  a random subset of  $Q$  of size  $n$  ;
 $T \leftarrow \sqrt{n}W_1(d_{\mathbb{1}_{P'},m}(\mathbb{1}_{P'}), d_{\mathbb{1}_{Q'},m}(\mathbb{1}_{Q'}))$  ;
# Compute  $boot$  a  $N_{MC}$ -sample from the bootstrap law
 $dtmP \leftarrow (d_{\mathbb{1}_P,m}(x))_{x \in P}$  ;
 $dtmQ \leftarrow (d_{\mathbb{1}_Q,m}(x))_{x \in Q}$  ;
Let  $boot$  be empty ;
for  $j$  in  $1 \dots \lfloor N_{MC}/2 \rfloor$  :
  Let  $dtmP_1$  and  $dtmP_2$  be two independent  $n$ -samples from  $\mathbb{1}_{dtmP}$  ;
  Let  $dtmQ_1$  and  $dtmQ_2$  be two independent  $n$ -samples from  $\mathbb{1}_{dtmQ}$  ;
  Add  $\sqrt{n}W_1(\mathbb{1}_{dtmP_1}, \mathbb{1}_{dtmP_2})$  and  $\sqrt{n}W_1(\mathbb{1}_{dtmQ_1}, \mathbb{1}_{dtmQ_2})$  to  $boot$  ;
# Compute  $qalph$ , the  $\alpha$ -quantile of  $boot$ 
Let  $qalph$  be the  $\lfloor N_{MC} - N_{MC} \times \alpha \rfloor$ th smallest element of  $boot$  ;
Output :  $(T \geq qalph)$ 

```

Recall that the L_1 -Wasserstein distance W_1 is simply the L_1 -norm of the difference between the cumulative distribution functions. It can be implemented by the **R** function `emd` from the package `emdist`. To compute the distance to an empirical measure at a point x , it is sufficient to search for its nearest neighbours; see section 2. This can be implemented by the **R** function `dtm` with tuning parameter $r = 1$, from the package `TDA` [15].

4.2 Validity of the method

In order to prove the validity of our method, we need to introduce a statistical framework.

First of all, from two N -samples from the mm-spaces $(\mathcal{X}, \delta, \mu)$ and $(\mathcal{Y}, \gamma, \nu)$, we derive four independent empirical measures, $\hat{\mu}_n, \hat{\mu}_{N-n}, \hat{\nu}_n$ and $\hat{\nu}_{N-n}$. We also denote $\hat{\mu}_N$ (respectively $\hat{\nu}_N$) the empirical measure associated to the whole N -sample of law μ (respectively ν), that is $\hat{\mu}_N = \frac{n}{N}\hat{\mu}_n + \frac{N-n}{N}\hat{\mu}_{N-n}$.

Then, we define the **test statistic** as:

$$T_{N,n,m}(\mu, \nu) = \sqrt{n}W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\hat{\nu}_N, m}(\hat{\nu}_n)).$$

Its law will be denoted by $\mathcal{L}_{N,n,m}(\mu, \nu)$.

Remark that for two isomorphic mm-spaces $(\mathcal{X}, \delta, \mu)$ and $(\mathcal{Y}, \gamma, \nu)$, the distribution of $T_{N,n,m}(\mu, \nu)$ is $\mathcal{L}_{N,n,m}(\mu, \mu)$, $\mathcal{L}_{N,n,m}(\nu, \nu)$, but also $\frac{1}{2}\mathcal{L}_{N,n,m}(\mu, \mu) + \frac{1}{2}\mathcal{L}_{N,n,m}(\nu, \nu)$; see Lemma 27 in the Appendix.

For some $\alpha > 0$, we denote by $q_\alpha = \inf\{x \in \mathbb{R} \mid F(x) \geq 1 - \alpha\}$, the α -**quantile** of a distribution with cumulative distribution function F .

The α -quantile $q_{\alpha, N, n}$ of $\frac{1}{2}\mathcal{L}_{N,n,m}(\mu, \mu) + \frac{1}{2}\mathcal{L}_{N,n,m}(\nu, \nu)$ will be approximated by the α -quantile $\hat{q}_{\alpha, N, n}$ of $\frac{1}{2}\mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N) + \frac{1}{2}\mathcal{L}_{N,n,m}^*(\hat{\nu}_N, \hat{\nu}_N)$. Here $\mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N)$ stands for the distribution of $\sqrt{n}W_1(d_{\hat{\mu}_N, m}(\mu_n^*), d_{\hat{\mu}_N, m}(\mu_n'^*))$ conditionally to $\hat{\mu}_N$, where μ_n^* and $\mu_n'^*$ are two empirical measures from independent n -samples of law $\hat{\mu}_N$.

The **test** we deal with in this paper is then:

$$\phi_N = \mathbb{1}_{T_{N,n,m}(\mu, \nu) \geq \hat{q}_{\alpha, N, n}}.$$

The null hypothesis H_0 is rejected if $\phi_N = 1$, that is if the L_1 -Wasserstein distance between the two empirical signatures $d_{\hat{\mu}_N, m}(\hat{\mu}_n)$ and $d_{\hat{\nu}_N, m}(\hat{\nu}_n)$ is too high.

4.2.1 A test of asymptotic level α

In this part, we prove that the test we propose is of asymptotic level α , that is such that:

$$\limsup_{N \rightarrow \infty} \mathbb{P}_{(\mu, \nu) \in H_0}(\phi_N = 1) \leq \alpha.$$

For this, we prove that the law of the test statistic $\frac{1}{2}\mathcal{L}_{N,n,m}(\mu, \mu) + \frac{1}{2}\mathcal{L}_{N,n,m}(\nu, \nu)$ under the hypothesis H_0 and the bootstrap law $\frac{1}{2}\mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N) + \frac{1}{2}\mathcal{L}_{N,n,m}^*(\hat{\nu}_N, \hat{\nu}_N)$ converge weakly to some fixed distribution when n and N go to ∞ . In order to adopt a non-asymptotic and more visual point of view, we also derive upper bounds in expectation for the L_1 -Wasserstein distance between these two distributions.

Remark that it is sufficient to prove weak convergence for $\mathcal{L}_{N,n,m}(\mu, \mu)$ and $\mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N)$. Moreover,

$$W_1\left(\frac{1}{2}\mathcal{L}_{N,n,m}(\mu, \mu) + \frac{1}{2}\mathcal{L}_{N,n,m}(\nu, \nu), \frac{1}{2}\mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N) + \frac{1}{2}\mathcal{L}_{N,n,m}^*(\hat{\nu}_N, \hat{\nu}_N)\right)$$

is upper bounded by

$$\frac{1}{2}W_1(\mathcal{L}_{N,n,m}(\mu, \mu), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N)) + \frac{1}{2}W_1(\mathcal{L}_{N,n,m}(\nu, \nu), \mathcal{L}_{N,n,m}^*(\hat{\nu}_N, \hat{\nu}_N)).$$

This is a straightforward consequence of the definition of the L_1 -Wasserstein distance with transport plans. Thus, this is also sufficient to derive upper bounds in expectation for the quantity $W_1(\mathcal{L}_{N,n,m}(\mu, \mu), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N))$.

Lemma 16.

For μ a measure supported on a compact set, we choose n as a function of N such that: when N goes to infinity, n goes to infinity, $\sqrt{n}\mathbb{E}[\|\mathbf{d}_{\mu,m} - \mathbf{d}_{\hat{\mu}_N,m}\|_{\infty,\mathcal{X}}]$ goes to zero or more specifically $\frac{\sqrt{n}}{m}\mathbb{E}[W_1(\mu, \hat{\mu}_N)]$ goes to zero. Then we have that:

$$\mathcal{L}_{N,n,m}(\mu, \mu) \rightsquigarrow \mathcal{L}(\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1),$$

when N goes to infinity. Moreover, if n is chosen such that $\sqrt{n}W_1(\mathbf{d}_{\mu,m}(\mu), \mathbf{d}_{\mu,m}(\hat{\mu}_N))$ and $\sqrt{n}\|\mathbf{d}_{\mu,m} - \mathbf{d}_{\hat{\mu}_N,m}\|_{\infty,\mathcal{X}}$ go to zero a.e., we have that for almost every sample $X_1, X_2, \dots, X_N, \dots$:

$$\mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N) \rightsquigarrow \mathcal{L}(\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1),$$

when N goes to infinity; with $\mathbb{G}_{\mu,m}$ and $\mathbb{G}'_{\mu,m}$ two independent Gaussian processes with covariance kernel $\kappa(s, t) = F_{\mathbf{d}_{\mu,m}(\mu)}(s)(1 - F_{\mathbf{d}_{\mu,m}(\mu)}(t))$ for $s \leq t$.

Proof

Proof in the Appendix, in Section C.3. ■

Proposition 17.

If the two weak convergences in lemma 16 occur, and if the α -quantile q_α of the distribution $\mathcal{L}(\frac{1}{2}\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1 + \frac{1}{2}\|\mathbb{G}_{\nu,m} - \mathbb{G}'_{\nu,m}\|_1)$ is a point of continuity of its cumulative distribution function, then the asymptotic level of the test at (μ, ν) is α .

Proof

Proof in the Appendix, in Section C.3. ■

Remark that for uniform measures on any sphere in \mathbb{R}^d , the continuity assumption for the cumulative distribution function of $\mathcal{L}(\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1)$ is not satisfied. This is a degenerated case. Thus, the test cannot be applied to such mm-spaces.

We choose $N = cn^\rho$ for some positive constants ρ and c . Then the test is asymptotically valid for two measures supported on a compact subset of the Euclidean space \mathbb{R}^d if we assume that $\rho > \frac{\max\{d, 2\}}{2}$.

Proposition 18.

Let μ be some Borel probability measure supported on some compact subset of \mathbb{R}^d . Under the assumption

$$\rho > \frac{\max\{d, 2\}}{2},$$

the two weak convergences of lemma 16 occur.

Moreover, a bound for the expectation of $W_1(\mathcal{L}_{N,n,m}(\mu, \mu), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N))$ is of order:

$$N^{\frac{1}{2\rho} - \frac{1}{\max\{d, 2\}}} (\log(1 + N))^{1_{d=2}}.$$

And, $W_1(\mathcal{L}_{N,n,m}(\mu, \mu), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N)) \rightarrow 0$ a.e. when n goes to ∞ .

Proof

This proposition is based on rates of convergence for the Wasserstein distance between a measure μ with values in \mathbb{R}^d and its empirical version $\hat{\mu}_N$; see [17] for general dimensions and [6] for $d = 1$. Proof in the Appendix, in Section C.4. ■

A probability measure μ is (a, b) -**standard** with positive parameters a and b , if for all positive radius r and any point x of the support of μ , we have that $\mu(\mathbb{B}(x, r)) \geq \min\{1, ar^b\}$. Uniform measures on open subsets of \mathbb{R}^d satisfy such a property:

Example 19. Let O be a non-empty bounded open subset of \mathbb{R}^d . Then, the measure μ_O is (a, d) -standard with

$$a = \frac{\omega_d}{\text{Leb}_d(O)} \left(\frac{\text{Reach}(O)}{\mathcal{D}(O)} \right)^d.$$

Here, $\mathcal{D}(O)$ stands for the diameter of O and ω_d for $\text{Leb}_d(\text{B}(0, 1))$, the Lebesgue volume of the unit d -dimensional ball.

Proof

Proof in the Appendix, in Section A.1. ■

Similar results can be obtained for uniform measures on compact submanifolds of dimension d . In [24] (lemma 5.3), the authors give a bound for a depending on the reach of the submanifold.

The test is asymptotically valid for two (a, b) -standard measures supported on compact connected subsets of \mathbb{R}^d if $\rho > 1$:

Proposition 20.

Let μ be an (a, b) -standard measure supported on a connected compact subset of \mathbb{R}^d . The two weak convergences of lemma 16 occur if the assumption $\rho > 1$ is satisfied. Moreover, a bound for the expectation of $W_1(\mathcal{L}_{N,n,m}(\mu, \mu), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N))$ is of order $N^{\frac{1}{2\rho} - \frac{1}{2}}$ up to a logarithm term.

Proof

This proposition is based on rates of convergence for the infinity norm between the distance to a measure and its empirical version; see [13]. Proof in the Appendix, in Section C.5. ■

Remark that we can achieve a rate close to the parametric rate for Ahlfors regular measures, whereas for general measures, the rate gets worse when the dimension increases. Anyway, we need ρ to be as big as possible for the bootstrapped law to be a good enough approximation of the law of the statistic, that is to have a type I error close enough to α ; keeping in mind that n should go to ∞ with N .

4.2.2 The power of the test

The **power** of the test $\phi_N = \mathbb{1}_{\sqrt{n}W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\hat{\nu}_N, m}(\hat{\nu}_n)) \geq \hat{q}_{\alpha, N, n}}$ is defined for two mm-spaces $(\mathcal{X}, \delta, \mu)$ and $(\mathcal{Y}, \gamma, \nu)$ by:

$$1 - \mathbb{P}_{(\mu, \nu)}(\phi_N = 0).$$

If the spaces are not isomorphic, we want the test to reject the null with high probability. It means that we want the power to be as big as possible. Here, we give a lower bound for the power, or more precisely an upper bound for $\mathbb{P}_{(\mu, \nu)}(\phi_N = 0)$, the **type II error**.

Proposition 21.

Let μ and ν be two Borel measures supported on \mathcal{X} and \mathcal{Y} , two compact subsets of \mathbb{R}^d . We assume that the mm-spaces $(\mathcal{X}, \delta, \mu)$ and $(\mathcal{Y}, \gamma, \nu)$ are non-isomorphic and that the DTM-signature is discriminative for some m in $(0, 1]$, that is such that $W_1(d_{\mu, m}(\mu), d_{\nu, m}(\nu)) > 0$. We choose $N = n^\rho$ with $\rho > 1$. Then for all positive ϵ , there exists n_0 depending on μ and ν such that for all $n \geq n_0$, the type II error

$$\mathbb{P}_{(\mu, \nu)}(\sqrt{n}W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\hat{\nu}_N, m}(\hat{\nu}_n)) < \hat{q}_{\alpha, N, n})$$

is upper bounded by

$$4 \exp \left(- \frac{W_1^2(d_{\mu, m}(\mu), d_{\nu, m}(\nu))}{(2 + \epsilon) \max\{\mathcal{D}_{\mu, m}^2, \mathcal{D}_{\nu, m}^2\}} n \right),$$

with $\mathcal{D}_{\mu,m}$, the diameter of the support of the measure $d_{\mu,m}(\mu)$.

Proof

Proof in the Appendix, in Section C.6. ■

In order to have a high power, that is to reject H_0 more often when the mm-spaces are not isomorphic, we need n to be big enough, that is ρ small enough. Recall that n has to be small enough for the law of the statistic and its bootstrap version to be close. It means that some compromise should be done. Moreover, the choice of m for the test should depend on the geometry of the mm-spaces. The tuning of these parameters from data is still an open question.

5 Numerical illustrations

Let μ_ν be the distribution of the random vector $(R \sin(vR) + 0.03N, R \cos(vR) + 0.03N')$ with R , N and N' independent random variables; N and N' from the standard normal distribution and R uniform on $(0, 1)$. With the notation given in the Introduction, we consider the sets $\Theta_0 = \{(\mu_{10}, \mu_{10})\}$ and $\Theta_1 = \{(\mu_{10}, \mu_p) \mid p \neq 10\}$. We sample $N = 2000$ points from two measure, choose $\alpha = 0.05$, $m = 0.05$, $n = 20$, and $N_{MC} = 1000$. We give an example under which our test (**DTM**) is working and more powerful than (**KS**), which consists in applying a Kolmogorov-Smirnov test to $\frac{N}{2}$ -samples from $\mathcal{L}(\delta(X, X'))$ and $\mathcal{L}(\gamma(Y, Y'))$ with X and X' (resp. Y and Y') independent from μ (resp. ν). The experiments are repeated 1000 times to approximate the type I error for our test and the power for both tests.

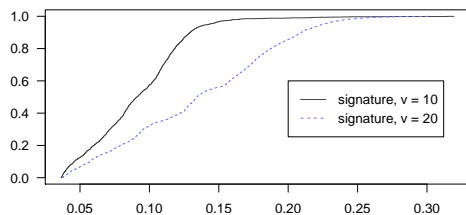


Figure 3: DTM-signature estimates, $m = 0.05$

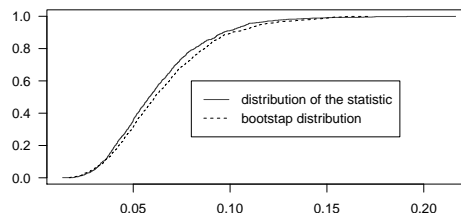


Figure 4: Bootstrap validity, $v = 10$, $m = 0.05$

v	15	20	30	40	100
type I error DTM	0.050	0.049	0.051	0.044	0.051
power DTM	0.525	0.884	0.987	0.977	0.985
power KS	0.768	0.402	0.465	0.414	0.422

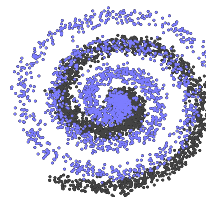


Figure 5: Type I error and power approximations

6 Concluding remarks and perspectives

This paper opens a new horizon of statistical tests based on shape signatures. It could be of interest to adapt these kind of methods to other signatures, if possible. In future it could even be interesting to build statistical tests based on many different signatures, leading to an even better discrimination. Regarding the test proposed in this paper itself, the geometric and statistical problem of the choice of the best parameters to use in practice is still an open, tough and engaging question.

Acknowledgements The author is extremely grateful to Frédéric Chazal, Pascal Massart and Bertrand Michel for introducing her to the distance to a measure, for their valuable comments and advises, and for proofreading.

References

- [1] Alejandro de Acosta and Evarist Giné. *Convergence Of Moments And Related Functionals In The Central Limit Theorem In Banach Spaces*. Z. Wahrsch.ver.Geb., 1979.
- [2] Aloisio Araujo and Evarist Giné. *The Central Limit Theorem for Real and Banach Valued Random Variables*. John Wiley & Sons Inc, 1980.
- [3] Eustasio del Barrio, Evarist Giné, and Carlos Matrán. “Central Limit Theorems For The Wasserstein Distance Between The Empirical And The True Distributions”. In: *The Annals of Probability* 27.2 (1999), pp. 1009–1071.
- [4] Peter J. Bickel and Kjell A. Doksum. *Mathematical statistics : basic ideas and selected topics*. Englewood Cliffs, N.J. Prentice Hall, 1977. ISBN: 0-13-564147-0. URL: <http://opac.inria.fr/record=b1089888>.
- [5] Patrick Billingsley. *Convergence of Probability Measures*. Wiley-Interscience, 1999.
- [6] Sergey Bobkov and Michel Ledoux. “One-Dimensional Empirical Measures, Order Statistics, And Kantorovich Transport Distances”. unpublished. 2014. URL: <http://perso.math.univ-toulouse.fr/ledoux/files/2014/04/Order.statistics.pdf>.
- [7] Mickaël Buchet. “Topological Inference From Measures”. PhD thesis. Université Paris-Sud – Paris XI, 2014.
- [8] Blanche Buet and Gian Paolo Leonardi. “Recovering Measures From Approximate Values On Balls”. unpublished. 2015. URL: <http://arxiv.org/abs/1510.02793>.
- [9] Frédéric Cazals and Alix Lhéritier. “Beyond Two-sample-tests: Localizing Data Discrepancies in High-dimensional Spaces”. In: *IEEE/ACM DSAA*. 2015.
- [10] F Chazal et al. “Gromov-Hausdorff Stable Signatures for Shapes using Persistence”. In: *Computer Graphics Forum (proc. SGP 2009)* (2009), pp. 1393–1403.
- [11] Frédéric Chazal, David Cohen-Steiner, and Quentin Mérigot. “Geometric Inference for Probability Measures”. In: *Foundations of Computational Mathematics* 11.6 (2011), pp. 733–751.
- [12] Frédéric Chazal, Vin De Silva, and Steve Oudot. “Persistence stability for geometric complexes”. In: *Geometriae Dedicata* 173.1 (2014), pp. 193–214.
- [13] Frédéric Chazal, Pascal Massart, and Bertrand Michel. “Rates Of Convergence For Robust Geometric Inference”. In: *Electronic Journal of Statistics* 10.2 (2016), pp. 2243–2286.
- [14] Eustasio Del Barrio, Hélène Lescornel, and Jean-Michel Loubes. “A statistical analysis of a deformation model with Wasserstein barycenters : estimation procedure and goodness of fit test”. unpublished. 2015. URL: <http://arxiv.org/pdf/1508.06465v2.pdf>.
- [15] Brittany Terese Fasy et al. “Introduction to the R package TDA”. In: *CoRR* abs/1411.1830 (2014). URL: <http://arxiv.org/abs/1411.1830>.
- [16] Herbert Federer. “Curvature Measures”. In: *Transactions of the American Mathematical Society* 93.3 (1959), pp. 418–491.

- [17] Nicolas Fournier and Arnaud Guillin. “On The Rate Of Convergence In Wasserstein Distance Of The Empirical Measure”. In: *Probability Theory & Related Fields* 162 (3 2015), pp. 707–738.
- [18] Arthur Gretton et al. “A Kernel Two-Sample Test”. In: *Journal of Machine Learning Research* 13 (2012), pp. 723–773.
- [19] Mikhail Gromov. *Metric Structures for Riemannian and Non-Riemannian Spaces*. Birkhäuser Basel, 2003.
- [20] André Lieutier. “Any Open Bounded Subset of \mathbb{R}^n Has the Same Homotopy Type Than Its Medial Axis”. In: *Computer Aided Geometric Design* 36.11 (2004), pp. 1029–1046.
- [21] Ulrike von Luxburg and Morteza Alamgir. “Density estimation from unweighted k-nearest neighbor graphs: a roadmap”. In: *NIPS*. 2013.
- [22] Pascal Massart. “The Tight Constant in the Dvoretzky-Kiefer-Wolfowitz Inequality”. In: *The Annals of Probability* 18.3 (1990), pp. 1269–1283.
- [23] Facundo Mémoli. “Gromov–Wasserstein Distances and the Metric Approach to Object Matching”. In: *Foundations of Computational Mathematics* 11.4 (2011), pp. 417–487.
- [24] Partha Niyogi, Steven Smale, and Shmuel Weinberger. “Finding the Homology of Submanifolds with High Confidence from Random Samples”. In: *Discrete and Computational Geometry* 39 (1 2008), pp. 419–441.
- [25] Robert Osada et al. “Shape Distributions”. In: *ACM Transactions on Graphics* 21 (4 2002), pp. 807–832.
- [26] Cédric Villani. *Topics in Optimal Transportation*. American Mathematical Society, 2003.

Appendix

A Uniform measures on open subsets of \mathbb{R}^d

In this part, we focus on some mm-spaces $(O, \|\cdot\|_2, \mu_O)$ where O stands for a non-empty bounded open subset of \mathbb{R}^d satisfying $(\overline{O})^\circ = O$. The measure μ_O , the medial axis $\mathcal{M}(O)$ and the reach $\text{Reach}(O)$ have been defined in Section 2. The object $\epsilon(m, O)$ is defined for some mass parameter m in $[0, 1]$ by

$$\epsilon(m, O) = \left(\frac{m \text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}}.$$

This is the radius of a ball included in O , with μ_O measure equal to m . For some positive ϵ , O_ϵ stands for the set of points in O which distance to ∂O is not smaller than ϵ :

$$O_\epsilon = \left\{ x \in O, \inf_{y \in \partial O} \|x - y\|_2 \geq \epsilon \right\}.$$

A.1 The distance to uniform measures

Here, we derive some properties of the spaces $(O, \|\cdot\|_2, \mu_O)$. We give a lower bound for the minimum of the distance to the measure μ_O and give a description of the points attaining this bound. Then, we use such considerations to prove identifiability of the measure μ_O from its distance-to-a-measure function. That is, to prove Proposition 7 of the paper.

First, we state some technical lemma proposed by Lieutier in [20].

Lemma 22.

If we define the *skeleton* $\text{Sk}(O)$ of the open set O as the set of centres of maximal balls (for the inclusion) included in O , then we get:

$$\mathcal{M}(O) \subset \text{Sk}(O) \subset \overline{\mathcal{M}(O)}.$$

Now we can formulate some technical lemma:

Lemma 23.

For any x in O , there exist a maximal ball for the inclusion, included in O and containing x .

Proof

Let us consider the class $\mathcal{S} = \{B(y, r) \mid r > 0 \text{ and } x \in B(y, r) \subset O\}$ of all non-empty open balls included in O and containing x . We are going to show that this class contains a maximal element by using the Zorn's lemma. For this, we need to show that the partially-ordered set \mathcal{S} is inductive, which means that any non-empty totally-ordered subclass \mathcal{T} of \mathcal{S} is upper bounded by some element of \mathcal{S} . Let \mathcal{T} be a non-empty totally-ordered subclass of \mathcal{S} . Set $R = \sup\{r > 0 \mid \exists y \in O, B(y, r) \in \mathcal{T}\}$ the supremum of the radii of all balls in \mathcal{T} . Since \mathcal{T} is non-empty and O is bounded, R is positive and finite. Let $(y_k)_{k \in \mathbb{N}}$ be a sequence of centres of balls in \mathcal{T} converging to a point y in \mathbb{R}^d such that the sequence of associated radii $(r_k)_{k \in \mathbb{N}}$ is non decreasing with R as a limit. Since \mathcal{T} is totally-ordered and the radii non decreasing, the union $\bigcup_{k \in \mathbb{N}} B(y_k, r_k)$ is non decreasing, equal to $B(y, R)$. Thus, $B(y, R)$ belongs to \mathcal{S} and upper bounds \mathcal{T} . So the class \mathcal{S} is inductive and thanks to the Zorn's lemma, it contains a maximal element. ■

Proof of Example 19: For any point x in O and $r > 0$, thanks to Lemma 23 there exist a maximal ball $B(x', r')$ included in $O \cap B(x, r)$ which contains x . Assume for the sake of contradiction that $r' < \min \left\{ \frac{r}{2}, \text{Reach}(O) \right\}$.

Since $r' < \frac{r}{2}$, the ball $\bar{B}(x', r')$ is included in $B(x, r)$ thus $B(x', r')$ is maximal in O . So x' belongs to $\text{Sk}(O)$, and thanks to Lemma 22, to $\overline{\mathcal{M}(O)}$. But $r' < \text{Reach}(O)$; this is absurd.

It follows that:

$$\mu_O(B(x, r)) \geq \mu_O \left(B \left(x', \min \left\{ \text{Reach}(O), \frac{r}{2} \right\} \right) \right).$$

So, for $r \leq 2\text{Reach}(O)$, since $2\text{Reach}(O) \leq \mathcal{D}(O)$ by considering a point on $\text{Sk}(O)$, we get:

$$\mu_O(B(x, r)) \geq r^d \left(\frac{\text{Reach}(O)}{\mathcal{D}(O)} \right)^d \frac{\omega_d}{\text{Leb}_d(O)},$$

which is also true for r in $[2\text{Reach}(O), \mathcal{D}(O)]$, whereas for $r \geq \mathcal{D}(O)$ we have $\mu_O(B(x, r)) = 1$. The choice of a in the lemma is thus relevant. ■

We now focus on the set of points in \mathbb{R}^d minimizing the distance to the measure μ_O . For this, we need some lemma.

Lemma 24.

If x in \mathbb{R}^d satisfies $\mu_O(B(x, \epsilon)) = \frac{\omega_d \epsilon^d}{\text{Leb}_d(O)}$, then $B(x, \epsilon) \subset O$.

Proof

If x in \mathbb{R}^d satisfies $\mu_O(B(x, \epsilon)) = \frac{\omega_d \epsilon^d}{\text{Leb}_d(O)}$, then, $\text{Leb}_d(O^c \cap B(x, \epsilon)) = 0$. Assume for the sake of contradiction that the set $O^c \cap B(x, \epsilon)$ is not empty. Since $(\bar{O})^\circ = O$, then the open subset $(O^c)^\circ \cap B(x, \epsilon)$ of $O^c \cap B(x, \epsilon)$ is not empty, thus of positive Lebesgue measure, which is absurd. So $B(x, \epsilon) \subset O$. ■

Proposition 25.

The constant $d_{\min} = \frac{d}{d+1} \left(\frac{m\text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}}$ is a lower bound for the distance to the measure μ_O over \mathbb{R}^d . Moreover, the set of points attaining this bound is exactly $O_{\epsilon(m, O)}$.

Proof

Remark that for all positive l smaller than m , we have:

$$\delta_{\mu, l}(x) \geq \left(\frac{l\text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}}.$$

Moreover, these inequalities are equalities for all points x in $O_{\epsilon(m, O)}$. By integrating, we get the lower bound d_{\min} for $x \mapsto d_{\mu, m}(x)$, and it is attained on $O_{\epsilon(m, O)}$.

Now take some point x in \mathbb{R}^d satisfying $d_{\mu, m}(x) = d_{\min}$. For almost all l smaller than m , we have: $\delta_{\mu, l}(x) = \left(\frac{l\text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}}$. In particular we get for these values of l that:

$$\mu \left(\bar{B} \left(x, \left(\frac{m\text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}} \right) \right) > l.$$

So, $\mu \left(B \left(x, \left(\frac{m\text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}} \right) \right) = m$, and thanks to Lemma 24, we get that $x \in O_{\epsilon(m, O)}$. ■

Proposition 26.

If $\text{Reach}(O) \geq \epsilon(m, O)$, then:

$$\{x \in \mathbb{R}^d \mid d_{\mu, m}(x) = d_{\min}\}^{\epsilon(m, O)} = O,$$

where for any set A , the notation A^ϵ stands for $\bigcup_{x \in A} \overline{B}(x, \epsilon)$, the ϵ -offset of A .

Proof

Remind thanks to Proposition 25 that $\{x \in \mathbb{R}^d \mid d_{\mu, m}(x) = d_{\min}\} = O_{\epsilon(m, O)}$. Moreover, $O_{\epsilon(m, O)}^{\epsilon(m, O)} \subset O$. Assume for the sake of contradiction that the set $O \setminus O_{\epsilon(m, O)}^{\epsilon(m, O)}$ is non-empty. Take a point x in this set and consider $B(x', r')$ a maximal ball containing x and included in O given by Lemma 23. Since $x \notin O_{\epsilon(m, O)}^{\epsilon(m, O)}$, we get that $r' < \epsilon(m, O)$. Moreover, x' belongs to $\text{Sk}(O)$ and so, thanks to Lemma 22, to $\overline{M}(O)$. Then, by continuity of the function distance to the compact set ∂O , $r' = d_{\partial O}(x') \geq \text{Reach}(O) \geq \epsilon(m, O)$, which is a contradiction. So, $O_{\epsilon(m, O)}^{\epsilon(m, O)} = O$. ■

A.2 The DTM-signature to discriminate between uniform and non uniform measures.

Proof of Proposition 14: As for Proposition 25, we get that for any point x in O :

$$d_{\mu_O, m}(x) \geq d_{\min} := \left(\frac{m \text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}} \frac{d}{1+d}.$$

We will lower bound the L_1 -Wasserstein distance between $d_{\mu_O, m}(\mu_O)$ and $d_{\nu, m}(\nu)$ by the integral of $F_{d_{\nu, m}(\nu)}$ over the interval $[0, d_{\min}]$, since $F_{d_{\mu_O, m}(\mu_O)}$ equals zero on this interval. We thus need to lower bound $F_{d_{\nu, m}(\nu)}(t)$ for all $t \leq d_{\min}$.

As for Proposition 25, for $\lambda \geq 1$, any point x of $\{f \geq \lambda\}_{\lambda^{-\frac{1}{d}} \left(\frac{m \text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}}}$ satisfies $d_{\nu, m}(x) \leq \frac{d_{\min}}{\lambda^{\frac{1}{d}}}$. Thus,

$$F_{d_{\nu, m}(\nu)} \left(\frac{d_{\min}}{\lambda^{\frac{1}{d}}} \right) \geq \nu \left(\{f \geq \lambda\}_{\lambda^{-\frac{1}{d}} \left(\frac{m \text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}}} \right).$$

And we get by denoting $\lambda(t)$ the real number λ satisfying $t = \frac{d_{\min}}{\lambda^{\frac{1}{d}}}$, that:

$$W_1(d_{\mu_O, m}(\mu_O), d_{\nu, m}(\nu)) \geq \int_{t=0}^{d_{\min}} \nu \left(\{f \geq \lambda(t)\}_{\lambda(t)^{-\frac{1}{d}} \left(\frac{m \text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}}} \right) dt.$$

Since a cumulative distribution function is non decreasing, we get:

$$\begin{aligned} W_1(d_{\mu_O, m}(\mu_O), d_{\nu, m}(\nu)) &\geq \\ &\int_{t=0}^{d_{\min}} \sup_{t' \leq t} \nu \left(\{f \geq \lambda(t')\}_{\lambda(t')^{-\frac{1}{d}} \left(\frac{m \text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}}} \right) dt \\ &= \int_{\lambda=1}^{\infty} d_{\min} \frac{1}{d} \frac{1}{\lambda^{\frac{1}{d}}} \frac{1}{\lambda} \sup_{\lambda' \geq \lambda} \nu \left(\{f \geq \lambda'\}_{\lambda'^{-\frac{1}{d}} \left(\frac{m \text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}}} \right) d\lambda \\ &\geq \frac{1}{d+1} \left(\frac{m \text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}} \int_{\lambda=1}^{\infty} \frac{1}{\lambda^{\frac{1}{d}}} \sup_{\lambda' \geq \lambda} \mu_O \left(\{f \geq \lambda'\}_{\left(\frac{m \text{Leb}_d(O)}{\lambda' \omega_d} \right)^{\frac{1}{d}}} \right) d\lambda. \end{aligned}$$

■

Now we assume that the density f is Hölder over O with parameters χ in $[0, 1]$ and L in \mathbb{R}_+^* .

Proof of Proposition 15: First remark that for all positive λ , with $\epsilon(\lambda) = \lambda^{-\frac{1}{d}} \left(\frac{m\text{Leb}_d(O)}{\omega_d} \right)^{\frac{1}{d}}$ we have:

$$\{f \geq \lambda + L\epsilon(\lambda)^x\} \cap O_{\epsilon(\lambda)} \subset \{f \geq \lambda\}_{\epsilon(\lambda)}.$$

According to Proposition 14, the aim is thus to show that for some λ bigger than 1, the set $\{f \geq \lambda + L\epsilon(\lambda)^x\} \cap O_{\epsilon(\lambda)}$ is non-empty. We thus focus on the supremum of f over $O_{\epsilon(\lambda)}$, which we denote by $\|f\|_{\infty, \epsilon(\lambda)}$.

Remind that if $\text{Reach}(O) \geq \epsilon(\lambda)$, then thanks to Proposition 26, the set $O_{\epsilon(\lambda)}^{\epsilon(\lambda)}$ equals O . Since f is Hölder, we can thus build some sequence $(y_n)_{n \in \mathbb{N}^*}$ in $O_{\epsilon(\lambda)}$, such that $f(y_n) \geq \|f\|_{\infty, O} - \frac{1}{n} - L\epsilon(\lambda)^x$. Finally we get:

$$\|f\|_{\infty, \epsilon(\lambda)} \geq \|f\|_{\infty, O} - L\epsilon(\lambda)^x.$$

So the quantity $W_1(d_{\mu_O, m}(\mu_O), d_{\nu, m}(\nu))$ is positive whenever:

$$\|f\|_{\infty, O} > \inf \{ \lambda + 2L\epsilon(\lambda)^x \mid \lambda \geq 1, \epsilon(\lambda) \leq \text{Reach}(O) \}.$$

With $\lambda_0 = 1$, we have $\lambda_0 + 2L\epsilon(\lambda_0)^x = 1 + 2L \left(\frac{m\text{Leb}_d(O)}{\omega_d} \right)^{\frac{x}{d}}$.

With λ_1 satisfying $\epsilon(\lambda_1) = \text{Reach}(O)$, we have:

$$\lambda_1 + 2L\epsilon(\lambda_1)^x = \frac{1}{(\text{Reach}(O))^d} \frac{m\text{Leb}_d(O)}{\omega_d} + 2L(\text{Reach}(O))^x.$$

We also have that

$$\inf \{ \lambda + 2L\epsilon(\lambda)^x \mid \lambda > 0 \} = (2L)^{\frac{d}{d+x}} \left(\frac{\text{Leb}_d(O)}{\omega_d} \right)^{\frac{x}{d+x}} m^{\frac{x}{d+x}} \left[\left(\frac{x}{d} \right)^{\frac{d}{d+x}} + \left(\frac{x}{d} \right)^{-\frac{x}{d+x}} \right].$$

The infimum is attained at $\lambda_2 = \left(\frac{x}{d} \right)^{\frac{d}{d+x}} (2L)^{\frac{d}{d+x}} \left(\frac{m\text{Leb}_d(O)}{\omega_d} \right)^{\frac{x}{d+x}}$.

It proves the first part of the proposition.

The second part is a straightforward consequence of the proof of Proposition 14. ■

B Stability of the DTM-signature

Proof of Proposition 9: The proof is relatively similar to the ones given by Mémoli in [23] for other signatures.

For any map plan π between μ and ν Borel measures on (\mathcal{X}, δ) and (\mathcal{Y}, γ) , we get:

$$\begin{aligned}
W_1(d_{\mu,m}(\mu), d_{\nu,m}(\nu)) &\leq \\
&\int_{\mathcal{X} \times \mathcal{Y}} |d_{\mu,m}(x) - d_{\nu,m}(y)| d\pi(x, y) = \\
&\int_{\mathcal{X} \times \mathcal{Y}} \left| \frac{1}{m} \int_0^m \delta_{\mu,l}(x) dl - \frac{1}{m} \int_0^m \delta_{\nu,l}(y) dl \right| d\pi(x, y) \leq \\
&\int_{\mathcal{X} \times \mathcal{Y}} \frac{1}{m} \int_0^m |\delta_{\mu,l}(x) - \delta_{\nu,l}(y)| dl d\pi(x, y) = \\
&\frac{1}{m} \int_{\mathcal{X} \times \mathcal{Y}} \int_0^m |\inf\{r > 0 \mid \mu(\overline{\mathbb{B}}(x, r)) > l\} - \inf\{r > 0 \mid \nu(\overline{\mathbb{B}}(y, r)) > l\}| dl d\pi(x, y) = \\
&\frac{1}{m} \int_{\mathcal{X} \times \mathcal{Y}} \int_0^m \left| \int_0^{+\infty} \left(\mathbb{1}_{\mu(\overline{\mathbb{B}}(x, r)) \leq l} - \mathbb{1}_{\nu(\overline{\mathbb{B}}(y, r)) \leq l} \right) dr \right| dl d\pi(x, y) \leq \\
&\frac{1}{m} \int_{\mathcal{X} \times \mathcal{Y}} \int_0^{+\infty} \int_0^m \left| \mathbb{1}_{\mu(\overline{\mathbb{B}}(x, r)) \leq l} - \mathbb{1}_{\nu(\overline{\mathbb{B}}(y, r)) \leq l} \right| dl dr d\pi(x, y) \leq \\
&\frac{1}{m} \int_{\mathcal{X} \times \mathcal{Y}} \int_0^{+\infty} |\mu(\overline{\mathbb{B}}(x, r)) \wedge m - \nu(\overline{\mathbb{B}}(y, r)) \wedge m| dr d\pi(x, y) \leq \\
&\frac{1}{m} \int_{\mathcal{X} \times \mathcal{Y}} \int_0^{+\infty} \left| \int_{\mathcal{X} \times \mathcal{Y}} \left(\mathbb{1}_{\delta(x, x') \leq r} - \mathbb{1}_{\gamma(y, y') \leq r} \right) d\pi(x', y') \right| \wedge m dr d\pi(x, y) \leq \\
&\frac{1}{m} \int_{\mathcal{X} \times \mathcal{Y}} \int_{\mathcal{X} \times \mathcal{Y}} \int_0^{+\infty} |\mathbb{1}_{\delta(x, x') \leq r} - \mathbb{1}_{\gamma(y, y') \leq r}| dr d\pi(x', y') d\pi(x, y) = \\
&\frac{1}{m} \int_{\mathcal{X} \times \mathcal{Y}} \int_{\mathcal{X} \times \mathcal{Y}} |\delta(x, x') - \gamma(y, y')| d\pi(x', y') d\pi(x, y),
\end{aligned}$$

which concludes.

C The test

C.1 A lemma

Lemma 27 (EQUALITY OF EMPIRICAL SIGNATURES UNDER THE ISOMORPHIC ASSUMPTION).

If $(\mathcal{X}, \delta, \mu)$ and $(\mathcal{Y}, \gamma, \nu)$ are two isomorphic mm-spaces, then the distributions of the random variables

$$\sqrt{n}W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\hat{\mu}'_N, m}(\hat{\mu}'_n))$$

and

$$\sqrt{n}W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\hat{\nu}_N, m}(\hat{\nu}_n))$$

are equal. Here the empirical measures are all independent and the measures $\hat{\mu}'_N$ and $\hat{\mu}'_n$ are from samples from μ .

Proof

Remark that for $(X'_1, X'_2, \dots, X'_N)$ a N -sample of law μ and ϕ an isomorphism between $(\mathcal{X}, \delta, \mu)$ and $(\mathcal{Y}, \gamma, \nu)$, the tuple $(\phi(X'_1), \phi(X'_2), \dots, \phi(X'_N))$ is a N -sample of law ν . Moreover, $\delta(X'_i, X'_j) = \gamma(\phi(X'_i), \phi(X'_j))$ for all i and j in $[[1, N]]$. It follows that the distances and the nearest neighbours are preserved.

Thus, the distributions of $(d_{\hat{\mu}_N, m}(X'_i))_{i \in [[1, n]]}$ and $(d_{\hat{\nu}_N, m}(Y_i))_{i \in [[1, n]]}$ are equal.

The lemma follows from the equality:

$$\begin{aligned} & W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\hat{\nu}_N, m}(\hat{\nu}_n)) \\ &= \int_0^{+\infty} \frac{1}{n} \left| \sum_{i=1}^n \mathbb{1}_{d_{\hat{\mu}_N, m}(X_i) \leq s} - \sum_{i=1}^n \mathbb{1}_{d_{\hat{\nu}_N, m}(Y_i) \leq s} \right| ds, \end{aligned}$$

with (X_1, X_2, \dots, X_N) a N -sample from μ . ■

C.2 L_1 -Wasserstein distance between the laws of interest

Lemma 28.

The quantity $W_1(\mathcal{L}_{N, n, m}(\mu, \mu), \mathcal{L}_{N, n, m}^*(\hat{\mu}_N, \hat{\mu}_N))$ is upper bounded by:

$$2\sqrt{n}(\mathbb{E}[\|d_{\hat{\mu}_N, m} - d_{\mu, m}\|_{\infty, \mathcal{X}}] + W_1(d_{\mu, m}(\mu), d_{\mu, m}(\hat{\mu}_N)) + \|d_{\mu, m} - d_{\hat{\mu}_N, m}\|_{\infty, \mathcal{X}}).$$

Proof

Let (X_1, X_2, \dots, X_N) be a N -sample of law μ , and $\hat{\mu}_N$ the associated empirical measure. We can upper bound the L_1 -Wasserstein distance between the bootstrap law $\mathcal{L}^*(\sqrt{n}W_1(d_{\hat{\mu}_N, m}(\mu_n^*), d_{\hat{\mu}_N, m}(\mu_n^*))|\hat{\mu}_N)$ and the law of interest $\mathcal{L}(\sqrt{n}W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\hat{\mu}'_N, m}(\hat{\mu}'_n)))$, by:

$$W_1(\mathcal{L}(\sqrt{n}W_1(d_{\hat{\mu}_N, m}(\mu_n^*), d_{\hat{\mu}_N, m}(\mu_n^*))|\hat{\mu}_N), \mathcal{L}(\sqrt{n}W_1(d_{\mu, m}(\mu_n^*), d_{\mu, m}(\mu_n^*))|\hat{\mu}_N)) \quad (1)$$

$$+ W_1(\mathcal{L}(\sqrt{n}W_1(d_{\mu, m}(\mu_n^*), d_{\mu, m}(\mu_n^*))|\hat{\mu}_N), \mathcal{L}(\sqrt{n}W_1(d_{\mu, m}(\hat{\mu}_n), d_{\mu, m}(\hat{\mu}'_n)))) \quad (2)$$

$$+ W_1(\mathcal{L}(\sqrt{n}W_1(d_{\mu, m}(\hat{\mu}_n), d_{\mu, m}(\hat{\mu}'_n))), \mathcal{L}(\sqrt{n}W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\hat{\mu}'_N, m}(\hat{\mu}'_n)))) \quad (3)$$

We bound the term 1 by:

$$2\sqrt{n}\|d_{\mu, m} - d_{\hat{\mu}_N, m}\|_{\infty, \mathcal{X}}.$$

the term 2 by

$$2\sqrt{n}W_1(d_{\mu, m}(\mu), d_{\mu, m}(\hat{\mu}_N))$$

and the term 3 by

$$2\sqrt{n}\mathbb{E}[\|d_{\mu, m} - d_{\hat{\mu}_N, m}\|_{\infty, \mathcal{X}}].$$

This is proved in the three following lemmata. ■

Lemma 29 (Study of term 3).

We have

$$\begin{aligned} & W_1(\mathcal{L}(\sqrt{n}W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\hat{\mu}'_N, m}(\hat{\mu}'_n))), \mathcal{L}(\sqrt{n}W_1(d_{\mu, m}(\hat{\mu}_n), d_{\mu, m}(\hat{\mu}'_n)))) \leq \\ & 2\sqrt{n}\mathbb{E}[\|d_{\mu, m} - d_{\hat{\mu}_N, m}\|_{\infty, \mathcal{X}}]. \end{aligned}$$

Proof

To bound this L_1 -Wasserstein distance, we choose as a transport plan the law of the random vector

$$(\sqrt{n}W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\hat{\mu}'_N, m}(\hat{\mu}'_n)), \sqrt{n}W_1(d_{\mu, m}(\hat{\mu}_n), d_{\mu, m}(\hat{\mu}'_n))),$$

with $\hat{\mu}_n$, $\hat{\mu}'_n$, $\hat{\mu}_{N-n}$ and $\hat{\mu}'_{N-n}$ independent empirical measures of law μ . Then the L_1 -Wasserstein distance is bounded by:

$$\mathbb{E}[\|\sqrt{n}W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\hat{\mu}'_N, m}(\hat{\mu}'_n)) - \sqrt{n}W_1(d_{\mu, m}(\hat{\mu}_n), d_{\mu, m}(\hat{\mu}'_n))\|],$$

which is not bigger than:

$$\sqrt{n}\mathbb{E}[W_1(d_{\hat{\mu}_N,m}(\hat{\mu}_n), d_{\mu,m}(\hat{\mu}_n)) + W_1(d_{\hat{\mu}'_N,m}(\hat{\mu}'_n), d_{\mu,m}(\hat{\mu}'_n))].$$

We bound the term $\mathbb{E}[W_1(d_{\hat{\mu}_N,m}(\hat{\mu}_n), d_{\mu,m}(\hat{\mu}_n))]$ by $\mathbb{E}[\|d_{\mu,m} - d_{\hat{\mu}_N,m}\|_{\infty,\mathcal{X}}]$, thanks to Lemma 32. ■

Lemma 30 (Study of term 2).

We have

$$W_1(\mathcal{L}(\sqrt{n}W_1(d_{\mu,m}(\hat{\mu}_n), d_{\mu,m}(\hat{\mu}'_n)), \mathcal{L}(\sqrt{n}W_1(d_{\mu,m}(\mu_n^*), d_{\mu,m}(\mu_n'^*))|\hat{\mu}_N)) \leq 2\sqrt{n}W_1(d_{\mu,m}(\mu), d_{\mu,m}(\hat{\mu}_N)).$$

Proof

Let π be the optimal transport plan associated to $W_1(d_{\mu,m}(\mu), d_{\mu,m}(\hat{\mu}_N))$; see the definition of the L_1 -Wasserstein with transport plans.

From a n -sample of law π , we get two empirical distributions $d_{\mu,m}(\hat{\mu}_n)$ and $d_{\mu,m}(\mu_n^*)$. Independently, from another n -sample of law π , we get $d_{\mu,m}(\hat{\mu}'_n)$ and $d_{\mu,m}(\mu_n'^*)$.

The L_1 -Wasserstein distance is then bounded by:

$$\sqrt{n}\mathbb{E}_{\pi^{\otimes n} \otimes \pi^{\otimes n}}[W_1(d_{\mu,m}(\hat{\mu}_n), d_{\mu,m}(\mu_n^*)) + W_1(d_{\mu,m}(\hat{\mu}'_n), d_{\mu,m}(\mu_n'^*))].$$

Now remark that, if we denote $\hat{\mu}_n = \sum_{i=1}^n \frac{1}{n} \delta_{Y_i}$ and $\mu_n^* = \sum_{i=1}^n \frac{1}{n} \delta_{Z_i}$, we have:

$$\begin{aligned} W_1(d_{\mu,m}(\hat{\mu}_n), d_{\mu,m}(\mu_n^*)) &= \int_{t=0}^{+\infty} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{d_{\mu,m}(Y_i) \leq t} - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{d_{\mu,m}(Z_i) \leq t} \right| dt \\ &\leq \frac{1}{n} \sum_{i=1}^n \int_{t=0}^{+\infty} |\mathbb{1}_{d_{\mu,m}(Y_i) \leq t} - \mathbb{1}_{d_{\mu,m}(Z_i) \leq t}| dt \\ &= \frac{1}{n} \sum_{i=1}^n |d_{\mu,m}(Y_i) - d_{\mu,m}(Z_i)|. \end{aligned}$$

So, the L_1 -Wasserstein distance is not bigger than

$$2\sqrt{n}\mathbb{E}[|d_{\mu,m}(Y) - d_{\mu,m}(Z)|],$$

with $(d_{\mu,m}(Y), d_{\mu,m}(Z))$ of law π , so we get the upper bound:

$$2\sqrt{n}(W_1(d_{\mu,m}(\mu), d_{\mu,m}(\hat{\mu}_N))).$$

■

Lemma 31 (Study of term 1).

We have

$$W_1(\mathcal{L}(\sqrt{n}W_1(d_{\mu,m}(\mu_n^*), d_{\mu,m}(\mu_n'^*))|\hat{\mu}_N), \mathcal{L}(\sqrt{n}W_1(d_{\hat{\mu}_N,m}(\mu_n^*), d_{\hat{\mu}_N,m}(\mu_n'^*))|\hat{\mu}_N)) \leq 2\sqrt{n}\|d_{\mu,m} - d_{\hat{\mu}_N,m}\|_{\infty,\mathcal{X}}.$$

Proof

It is the same proof as for the first lemma, except that $\hat{\mu}_N$ is fixed. ■

Lemma 32.

Let ν , μ and μ' be some measures over some metric space (\mathcal{X}, δ) , we have:

$$W_1(d_{\mu,m}(\nu), d_{\mu',m}(\nu)) \leq \int_{\mathcal{X}} |d_{\mu,m}(x) - d_{\mu',m}(x)| d\nu(x) \leq \|d_{\mu,m} - d_{\mu',m}\|_{\infty, \text{Supp}(\nu)}.$$

Proof

We chose the transport plan $(d_{\mu,m}(Y), d_{\mu',m}(Y))$ for Y of law ν . ■

Thanks to Proposition 5 and to the fact that the distance to a measure is 1-Lipschitz, we can derive another upper bound depending only on the L_1 -Wasserstein distance between the measure μ and its empirical versions:

Corollary 33.

The quantity $W_1(\mathcal{L}_{N,n,m}(\mu, \mu), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N))$ is upper bounded by:

$$2\frac{\sqrt{n}}{m}\mathbb{E}[W_1(\hat{\mu}_N, \mu)] + 2\sqrt{n}\left(1 + \frac{1}{m}\right)W_1(\hat{\mu}_N, \mu).$$

The rates of convergence of the L_1 -Wasserstein distance between a Borel probability measure on the Euclidean space \mathbb{R}^d and its empirical version are faster when the dimension d is low; see [17]. Thus, we prefer to use the first bound for regular measures. In this case, we use rates of convergence for the distance to a measure, derived in [13]. For regular measures, in some cases, the bound in Lemma 28 is better than the bound in Corollary 33.

C.3 An asymptotic result with the convergence to the law of

$$\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1$$

Proof of Lemma 16: The random function $\sqrt{n}(F_{d_{\mu,m}(\mu)} - F_{d_{\mu,m}(\hat{\mu}_n)})$ converges weakly in L_1 to some gaussian process $\mathbb{G}_{\mu,m}$ with covariance kernel $\kappa(s, t) = F_{d_{\mu,m}(\mu)}(s)(1 - F_{d_{\mu,m}(\mu)}(t))$ for $s \leq t$; see [3] or part 3.3 of [6]. Thanks to Theorem 2.8 in [5], since $L_1 \times L_1$ is separable and $\hat{\mu}_n$ and $\hat{\mu}'_n$ are independent, the random vector

$$(\sqrt{n}(F_{d_{\mu,m}(\mu)} - F_{d_{\mu,m}(\hat{\mu}_n)}), \sqrt{n}(F_{d_{\mu,m}(\mu)} - F_{d_{\mu,m}(\hat{\mu}'_n)}))$$

converges weakly to $(\mathbb{G}_{\mu,m}, \mathbb{G}'_{\mu,m})$ with $\mathbb{G}_{\mu,m}$ and $\mathbb{G}'_{\mu,m}$ independent Gaussian processes. Since the map $(x, y) \mapsto x - y$ is continuous in L_1 , the mapping theorem states that $\sqrt{n}(F_{d_{\mu,m}(\hat{\mu}'_n)} - F_{d_{\mu,m}(\hat{\mu}_n)})$ converges weakly to the Gaussian process $\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}$ in L_1 . Once more we use the mapping theorem with the continuous map $x \mapsto \|x\|_1$ and the definition of the L_1 -Wasserstein distance as the L_1 -norm of the cumulative distribution functions to get that:

$$\sqrt{n}W_1(d_{\mu,m}(\hat{\mu}_n), d_{\mu,m}(\hat{\mu}'_n)) \rightsquigarrow \|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1.$$

We then get the convergence of moments following the same method as for Theorem 2.4 in [3]. We have the bound $\mathbb{E}[\|t \mapsto \mathbb{1}_{d_{\mu,m}(X_i) \leq t} - \mathbb{1}_{d_{\mu,m}(Y_i) \leq t}\|_1] \leq \mathcal{D}_\mu < \infty$. Moreover, the random function $\sqrt{n}(F_{d_{\mu,m}(\hat{\mu}'_n)} - F_{d_{\mu,m}(\hat{\mu}_n)})$ converges weakly to the gaussian process $\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}$ in L_1 . So, thanks to Theorem 5.1 in [1] (cited in [2] p.136), we have:

$$\mathbb{E}[\sqrt{n}W_1(d_{\mu,m}(\hat{\mu}_n), d_{\mu,m}(\hat{\mu}'_n))] \rightarrow \mathbb{E}[\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1].$$

We deduce that:

$$W_1(\mathcal{L}(\sqrt{n}W_1(d_{\mu,m}(\hat{\mu}_n), d_{\mu,m}(\hat{\mu}'_n))), \mathcal{L}(\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1)) \rightarrow 0.$$

Moreover, we have the bound:

$$W_1(\mathcal{L}(\sqrt{n}W_1(d_{\mu,m}(\hat{\mu}_n), d_{\mu,m}(\hat{\mu}'_n))), \mathcal{L}_{N,n,m}(\mu, \mu)) \leq 2\sqrt{n}\mathbb{E}[\|d_{\mu,m} - d_{\hat{\mu}_N,m}\|_{\infty, \mathcal{X}}].$$

So, if $\sqrt{n}\mathbb{E}[\|d_{\mu,m} - d_{\hat{\mu}_N,m}\|_{\infty, \mathcal{X}}] \rightarrow 0$ when $N \rightarrow \infty$, we have that:

$$W_1(\mathcal{L}_{N,n,m}(\mu, \mu), \mathcal{L}(\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1)) \rightarrow 0.$$

Finally, with the same arguments as for Lemma 28, we get that:

$$\begin{aligned} & W_1(\mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N), \mathcal{L}(\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1)) \leq \\ & W_1(\mathcal{L}(\sqrt{n}W_1(d_{\mu,m}(\hat{\mu}_n), d_{\mu,m}(\hat{\mu}'_n)), \mathcal{L}(\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1)) \\ & + 2\sqrt{n}W_1(d_{\mu,m}(\mu), d_{\mu,m}(\hat{\mu}_N)) + 2\sqrt{n}\|\mathbb{d}_{\mu,m} - \mathbb{d}_{\hat{\mu}_N,m}\|_{\infty, \mathcal{X}}). \end{aligned}$$

■

Proof of Proposition 17: Let $\epsilon < \alpha$ and η be two positive numbers.

The probability $\mathbb{P}_{(\mu,\nu)}(\phi_N = 1)$ is upper bounded by

$$\mathbb{P}(\sqrt{n}W_1(d_{\hat{\mu}_N,m}(\hat{\mu}_n), d_{\hat{\nu}_N,m}(\hat{\nu}_n)) \geq \mathfrak{q}_{\alpha+\epsilon} - \eta) + \mathbb{P}(\hat{\mathfrak{q}}_\alpha < \mathfrak{q}_{\alpha+\epsilon} - \eta).$$

With a drawing, we see that $\mathbb{P}(\hat{\mathfrak{q}}_\alpha < \mathfrak{q}_{\alpha+\epsilon} - \eta)$ is upper bounded by

$$\mathbb{P}\left(W_1\left(\mathcal{L}\left(\frac{1}{2}\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1 + \frac{1}{2}\|\mathbb{G}_{\nu,m} - \mathbb{G}'_{\nu,m}\|_1\right), \mathcal{L}^*\right) \geq \epsilon\eta\right),$$

where $\mathcal{L}^* = \frac{1}{2}\mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N) + \frac{1}{2}\mathcal{L}_{N,n,m}^*(\hat{\nu}_N, \hat{\nu}_N)$.

Thanks to the weak convergences in Lemma 16 of the paper and the Portmanteau lemma, $\limsup_{N \rightarrow \infty} \mathbb{P}_{(\mu,\nu)}(\phi_N = 1)$ is thus upper bounded by

$$\mathbb{P}\left(\frac{1}{2}\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1 + \frac{1}{2}\|\mathbb{G}_{\nu,m} - \mathbb{G}'_{\nu,m}\|_1 \geq \mathfrak{q}_{\alpha+\epsilon} - \eta\right).$$

We now make η and ϵ go to zero and under the continuity assumption, $\limsup_{N \rightarrow \infty} \mathbb{P}_{(\mu,\nu)}(\phi_N = 1) \leq \alpha$.

As well, we get that $\liminf_{N \rightarrow \infty} \mathbb{P}_{(\mu,\nu)}(\phi_N = 1) \geq \alpha$.

C.4 The case of measures supported on a compact subset of \mathbb{R}^d

Proof of part 2 of Proposition 18: We may assume that the diameter \mathcal{D}_μ of the support of the measure μ equals 1. Indeed, if we apply a dilatation to the measure to make the diameter of its support be equal to 1, then the quantity $W_1(\mathcal{L}_{N,n,m}(\mu, \mu), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N))$ is simply multiplied by the parameter of the dilatation. By using Corollary 33 and Theorem 1 of [17], we have a bound for the expectation:

$$\mathbb{E}\left[W_1(\mathcal{L}_{N,n,m}(\mu, \mu), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N))\right] \leq \begin{cases} C \frac{\sqrt{n}}{m} N^{-\frac{1}{d}} & \text{if } d > 2 \\ C \frac{\sqrt{n}}{m} N^{-\frac{1}{2}} \log(1+N) & \text{if } d = 2 \\ C \frac{\sqrt{n}}{m} N^{-\frac{1}{2}} & \text{if } d < 2 \end{cases}$$

for some positive constant C depending on μ . ■

Proof of part 3 of Proposition 18: First remark that for $\lambda > 1$,

$$\mathbb{P}(W_1(\mathcal{L}_{N,n,m}(\mu, \mu), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N)) \geq \lambda) = 0$$

under the assumption $\mathcal{D}_\mu = 1$. We thus focus on values of λ not bigger than 1. In this case, with the Theorem 2 of [17], we get easily that:

$$\mathbb{P} (W_1 (\mathcal{L}_{N,n,m}(\mu, \mu), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N)) \geq \lambda) \leq \begin{cases} C \exp \left(-C' \left(\lambda \frac{N^{\frac{1}{d}} m}{\sqrt{n}} - C'' \right)^d \right) & \text{for } d > 2 \\ C \exp \left(-C' \left(\frac{\frac{\sqrt{Nm}}{\sqrt{n}} \lambda - C'' \sqrt{\frac{N}{N-n}} \log(1+N-n)}{\log \left(2 + \frac{2\sqrt{N}}{\frac{\sqrt{Nm}}{\sqrt{n}} \lambda - C'' \sqrt{\frac{N}{N-n}} \log(1+N-n)} \right)} \right)^2 \right) & \text{for } d = 2 \\ C \exp \left(-C' \left(\lambda \frac{\sqrt{Nm}}{\sqrt{n}} - C'' \right)^2 \right) & \text{for } d < 2 \end{cases}$$

for some positive constants C , C' and C'' depending on μ .
We conclude the proof with the Borel–Cantelli lemma. ■

Proof of part 1 of Proposition 18: We need to show that under the assumption $\rho > \frac{\max\{d,2\}}{2}$, the following properties are satisfied:

$$\begin{aligned} \sqrt{n} \mathbb{E}[\|d_{\mu,m} - d_{\hat{\mu}_N,m}\|_{\infty, \mathcal{X}}] &\rightarrow 0, \\ \sqrt{n} W_1(d_{\mu,m}(\mu), d_{\mu,m}(\hat{\mu}_N)) &\rightarrow 0 \text{ a.e.}, \end{aligned}$$

and

$$\sqrt{n} \|d_{\mu,m} - d_{\hat{\mu}_N,m}\|_{\infty, \mathcal{X}} \rightarrow 0 \text{ a.e.}$$

We treat the case $d > 2$. The cases $d < 2$ and $d = 2$ are similar.

Thanks to Theorem 1 of [17], there is some positive constant C depending on μ such that for N big enough:

$$\mathbb{E}[W_1(\hat{\mu}_N, \mu)] \leq CN^{-\frac{1}{d}}.$$

Thus, thanks to part 2 of Proposition 18, the quantity $\sqrt{n} \mathbb{E}[\|d_{\mu,m} - d_{\hat{\mu}_N,m}\|_{\infty, \mathcal{X}}]$ goes to zero if $\frac{\sqrt{n}}{m} N^{-\frac{1}{d}}$ goes to zero when N goes to infinity. So, this convergence occurs under the assumption $\rho > \frac{d}{2}$.

We get from Theorem 2 of [17] that for $x \leq 1$, there are some positive constants C and c depending on μ such that:

$$\mathbb{P}(W_1(\hat{\mu}_N, \mu) \geq x) \leq C \exp(-cNx^d).$$

We use this inequality with $x = \frac{m}{\sqrt{n}} \frac{1}{K}$ for positive integers K . Thanks to the Borel–Cantelli lemma, under the assumption $\rho > \frac{d}{2}$, we get that:

$$\frac{\sqrt{n}}{m} W_1(\mu, \hat{\mu}_N) \rightarrow 0 \text{ a.e.}$$

So, thanks to Proposition 5, the third property is true.

To finish, remark that $d_{\mu,m}(\hat{\mu}_N)$ is the empirical measure associated to $d_{\mu,m}(\mu)$. Once more we use Theorem 2 of [17] and get that for $x \leq 1$, $\mathbb{P}(\sqrt{n} W_1(d_{\mu,m}(\hat{\mu}_N), d_{\mu,m}(\mu)) \geq x) \leq C \exp(-c \frac{N}{n} x^2)$. Thanks to the Borel–Cantelli lemma, under the assumption $\rho > 1$, the a.e. convergence to zero of $\sqrt{n} W_1(d_{\mu,m}(\mu), d_{\mu,m}(\hat{\mu}_N))$ occurs. ■

C.5 The case of (a, b) -standard measures

Let μ be a Borel probability measure supported on a connected compact subset \mathcal{X} of \mathbb{R}^d . We assume this measure to be (a, b) -standard for some positive numbers a and b . In this part, we derive rates of convergence in probability and in expectation for the quantity $\|\mathrm{d}_{\hat{\mu}_N, m} - \mathrm{d}_{\mu, m}\|_{\infty, \mathcal{X}}$. Thanks to these results, we can derive upper bounds and rates of convergence in expectation for $W_1(\mathcal{L}_{N, n, m}(\mu, \mu), \mathcal{L}_{N, n, m}^*(\hat{\mu}_N, \hat{\mu}'_N))$. We finally propose a choice for the parameter N depending on n for which the weak convergences $\mathcal{L}_{N, n, m}(\mu, \mu) \rightsquigarrow \|\mathbb{G}_{\mu, m} - \mathbb{G}'_{\mu, m}\|_1$ and $\mathcal{L}_{N, n, m}^*(\hat{\mu}_N, \hat{\mu}'_N) \rightsquigarrow \|\mathbb{G}_{\mu, m} - \mathbb{G}'_{\mu, m}\|_1$ occur.

C.5.1 Upper bounds for $\mathbb{P}(\sqrt{n}\|\mathrm{d}_{\hat{\mu}_N, m} - \mathrm{d}_{\mu, m}\|_{\infty, \mathcal{X}} \geq \lambda)$

We use the bounds given in Theorem 1 of [13], with the bound for the modulus of continuity given by **Lemma 3** in [13]: $\omega(h) = (\frac{h}{a})^{\frac{1}{b}}$. We directly get the following lemma:

Lemma 34 (Upper bound for $|\mathrm{d}_{\hat{\mu}_N, m}(x) - \mathrm{d}_{\mu, m}(x)|$).

Let x be a fixed point in \mathcal{X} and λ a positive number. We have,

$$\frac{1}{2}\mathbb{P}(|\mathrm{d}_{\hat{\mu}_N, m}(x) - \mathrm{d}_{\mu, m}(x)| \geq \lambda) \leq \exp\left(-2a^{\frac{2}{b}}Nm^{\frac{2b-2}{b}}\lambda^2\right) + \exp\left(-\frac{a}{2^{b-1}}N^{\frac{b+1}{2}}m^b\lambda^b\right) + \exp\left(-a^{\frac{1}{b}}N^{\frac{b+1}{2b}}m\lambda\right).$$

In order to derive an upper bound for $\|\mathrm{d}_{\hat{\mu}_N, m} - \mathrm{d}_{\mu, m}\|_{\infty, \mathcal{X}}$, like in [13], we use the fact that the function distance to a measure is 1-Lipschitz and that \mathcal{X} is compact, which means that we can compute a bound by upper-bounding the difference $|\mathrm{d}_{\hat{\mu}_N, m}(x) - \mathrm{d}_{\mu, m}(x)|$ over a finite number of points x of \mathcal{X} . Thanks to the following lemma, the minimal number of points needed for this purpose is not bigger than $\frac{(4\mathcal{D}_\mu\sqrt{d}+\lambda)^d}{\lambda^d}$:

Lemma 35.

Let μ is a measure supported on \mathcal{X} a compact subset of \mathbb{R}^d , and for $\lambda > 0$ denote $N(\mu, \lambda) = \inf\{N \in \mathbb{N}, \exists x_1, x_2 \dots x_N \in \mathcal{X}, \bigcup_{i \in [1, N]} \mathbb{B}(x_i, \lambda) \supset \mathcal{X}\}$. Then, we have:

$$N(\mu, \lambda) \leq \frac{(\mathcal{D}_\mu\sqrt{d} + \lambda)^d}{\lambda^d}.$$

Proof

The idea is to put a grid on the hypercube containing \mathcal{X} with edges of length \mathcal{D}_μ . The grid is a union of small hypercubes with edges of length equal to $\frac{\lambda}{\sqrt{d}}$, so that the number of such small hypercubes into which the big one is split is not superior to $(\frac{\mathcal{D}_\mu\sqrt{d}}{\lambda} + 1)^d$.

Then, we decide that each time the intersection between \mathcal{X} and some small hypercube is non-empty, we keep one of the elements of the intersection. We denote x_i the element associated to the i -th hypercube. Finally, each point x in \mathcal{X} belongs to a small hypercube, and its distance to the corresponding x_i is smaller than $\sqrt{\sum_{k=1}^d \frac{\lambda^2}{d}} = \lambda$. ■

We thus derive upper bounds for $\sqrt{n}\|\mathrm{d}_{\hat{\mu}_N, m} - \mathrm{d}_{\mu, m}\|_{\infty, \mathcal{X}}$:

Proposition 36 (Upper bound for $\sqrt{n}\|\mathrm{d}_{\hat{\mu}_N, m} - \mathrm{d}_{\mu, m}\|_{\infty, \mathcal{X}}$).

We have,

$$\begin{aligned} & \frac{\lambda^d}{2(4\mathcal{D}_\mu\sqrt{d} + \lambda)^d} \mathbb{P}(\sqrt{n}\|d_{\hat{\mu}_N, m} - d_{\mu, m}\|_{\infty, \mathcal{X}} \geq \lambda) \leq \\ & \exp\left(-\frac{a^{\frac{2}{b}} Nm^{\frac{2b-2}{b}}}{2n}\lambda^2\right) + \exp\left(-\frac{a}{2^{2b-1}} \frac{N^{\frac{b+1}{2}} m^b}{n^{\frac{b}{2}}}\lambda^b\right) + \exp\left(-\frac{a^{\frac{1}{b}} N^{\frac{b+1}{2b}} m}{2n^{\frac{1}{2}}}\lambda\right). \end{aligned}$$

Proof

Since the function distance to a measure is 1-Lipschitz, we get that:

$$\|d_{\hat{\mu}_N, m} - d_{\mu, m}\|_{\infty, \mathcal{X}} \leq \frac{\lambda}{2} + \sup_i \{|d_{\hat{\mu}_N, m}(x_i) - d_{\mu, m}(x_i)|\},$$

for the family $(x_i)_i$ associated to a grid which sides are of length equal to $\frac{\lambda}{4\sqrt{d}}$. We can thus bound the probability $\mathbb{P}(\|d_{\hat{\mu}_N, m} - d_{\mu, m}\|_{\infty, \mathcal{X}} \geq \lambda)$ by:

$$\sum_{i=1}^{N(\mu, \frac{\lambda}{4})} \mathbb{P}\left(|d_{\hat{\mu}_N, m}(x_i) - d_{\mu, m}(x_i)| \geq \frac{\lambda}{2}\right),$$

with $N(\mu, \frac{\lambda}{4}) \leq \frac{(4\mathcal{D}_\mu\sqrt{d} + \lambda)^d}{\lambda^d}$ thanks to Lemma 35. ■

C.5.2 Upper bounds for the expectation $\mathbb{E}[\|d_{\hat{\mu}_N, m} - d_{\mu, m}\|_{\infty, \mathcal{X}}]$

In order to get upper bounds for $\mathbb{E}[\|d_{\hat{\mu}_N, m} - d_{\mu, m}\|_{\infty, \mathcal{X}}]$, we use the same trick as used in [13], which is:

Lemma 37.

Let X a random variable such that:

$$\mathbb{P}(X \geq \lambda) \leq 1 \wedge D\lambda^{-q} \exp(-c\lambda^s)$$

for some integers q and s and some $D > 0$.

We have:

$$\mathbb{E}[X] \leq \left(\frac{\ln c}{c}\right)^{\frac{1}{s}} \left(\frac{q}{s}\right)^{\frac{1}{s}} \left[1 + D \left(\frac{q}{s}\right)^{\frac{-q-s}{s}} \frac{(\ln c)^{\frac{-q-s}{s}}}{s}\right].$$

More particularly, if $c \geq \exp D^{\frac{s}{q+s}} \frac{s}{q}$, then:

$$\mathbb{E}[X] \leq 2 \left(\frac{\ln c}{c}\right)^{\frac{1}{s}} \left(\frac{q}{s}\right)^{\frac{1}{s}}.$$

Proof

For any $\lambda_0 > 0$, that we can choose as $\lambda_0 = \frac{[\ln K]^{\frac{1}{s}}}{c^{\frac{1}{s}}}$, we get that:

$$\begin{aligned}
\mathbb{E}[X] &\leq \lambda_0 + \int_{\lambda_0}^{\infty} D\lambda^{-q} \exp(-c\lambda^s) d\lambda \\
&\leq \lambda_0 + D \frac{\lambda_0^{-q-s+1}}{cs} \exp -c\lambda_0^s \\
&= \frac{[\ln K]^{\frac{1}{s}}}{c^{\frac{1}{s}}} + D \frac{[\ln K]^{-\frac{q-s+1}{s}}}{sc c^{-\frac{q-s+1}{s}}} \frac{1}{K} \\
&= \frac{[\ln K]^{\frac{1}{s}}}{c^{\frac{1}{s}}} + \frac{[\ln K]^{\frac{1}{s}}}{c^{\frac{1}{s}}} D \frac{[\ln K]^{-\frac{q-s}{s}}}{sc^{\frac{-q}{s}}} \frac{1}{K} \\
&= \frac{[\ln K]^{\frac{1}{s}}}{c^{\frac{1}{s}}} \left[1 + D \frac{[\ln K]^{-\frac{q-s}{s}}}{sK c^{\frac{-q}{s}}} \right]
\end{aligned}$$

Finally, if we choose $K = c^{\frac{q}{s}}$, we get:

$$\mathbb{E}[X] \leq \left(\frac{q}{s}\right)^{\frac{1}{s}} \left[\frac{\ln c}{c}\right]^{\frac{1}{s}} \left[1 + D \left[\frac{q}{s}\right]^{\frac{-q-s}{s}} \frac{(\ln c)^{-\frac{q-s}{s}}}{s} \right].$$

■

From this lemma, we can derive the following lemma.

Lemma 38.

We have,

$$\begin{aligned}
&\mathbb{E}[\sqrt{n} \|d_{\hat{\mu}_{N,m}} - d_{\mu,m}\|_{\infty, \mathcal{X}}] \leq \\
&\square'_1 \frac{n^{\frac{1}{2}}}{N^{\frac{1}{2}} m^{\frac{b-1}{b}}} \left(\log \left(\frac{Nm^{\frac{2b-2}{b}}}{n} \right) \right)^{\frac{1}{2}} + \\
&\square'_2 \frac{n^{\frac{1}{2}}}{N^{\frac{b+1}{2b}} m} \left(\log \left(\frac{N^{\frac{b+1}{2}} m^b}{n^{\frac{b}{2}}} \right) \right)^{\frac{1}{b}} + \\
&\square'_3 \frac{n^{\frac{1}{2}}}{N^{\frac{b+1}{2b}} m} \log \left(\frac{N^{\frac{b+1}{2b}} m}{n^{\frac{1}{2}}} \right).
\end{aligned}$$

for some constants \square depending on a and b .

C.5.3 Upper bounds for the expectation of $W_1(\mathcal{L}_{N,n,m}(\mu, \mu), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N))$

Proof of part 2 of Proposition 20: For all $\lambda > 0$, for any measure μ Ahlfors b -regular with parameters (a, ∞) supported on a connected compact subset of \mathbb{R}^d , we can use Lemma 28 and Lemma 38 together with the rates of convergence of the L_1 -Wasserstein distance between empirical and true distribution in [6] to get the following result.

If $m \geq \frac{1}{2}$, then for n big enough we have, for some constants \square depending on a and b :

$$\begin{aligned} & \mathbb{E} \left[W_1 \left(\mathcal{L}_{N,n,m}(\mu, \mu), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N) \right) \right] \leq \\ & \square'_1 \frac{n^{\frac{1}{2}}}{(N)^{\frac{1}{2}} m^{\frac{b-1}{b}}} \left(\log \left(\frac{Nm^{\frac{2b-2}{b}}}{n} \right) \right)^{\frac{1}{2}} \\ & + \square'_2 \frac{n^{\frac{1}{2}}}{(N)^{\frac{b+1}{2b}} m} \left(\log \left(\frac{N^{\frac{b+1}{2}} m^b}{n^{\frac{b}{2}}} \right) \right)^{\frac{1}{b}} \\ & + \square'_3 \frac{n^{\frac{1}{2}}}{(N)^{\frac{b+1}{2b}} m} \log \left(\frac{N^{\frac{b+1}{2b}} m}{n^{\frac{1}{2}}} \right) \\ & + \square'_4 \frac{n^{\frac{1}{2}}}{N^{\frac{1}{2}}}. \end{aligned}$$

C.5.4 Convergence to the law of $\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1$

Proof of part 1 of Proposition 20: In order to get these two results, we use Lemma 16. The convergence to zero of $\sqrt{n}\mathbb{E}[\|d_{\mu,m} - d_{\hat{\mu}_{N-n},m}\|_{\infty,\mathcal{X}}]$ is a direct consequence of Lemma 38. We can derive a bound of its rate of convergence in $n^{\frac{1}{2}-\frac{\rho}{2}}$, up to a logarithm term. The a.e. convergence of $\sqrt{n}W_1(d_{\mu,m}(\mu), d_{\mu,m}(\hat{\mu}_N))$ to zero is derived as in the proof of Proposition 18, with the assumption $\rho > 1$. Finally, the a.e. convergence of $\sqrt{n}\|d_{\mu,m} - d_{\hat{\mu}_N,m}\|_{\infty,\mathcal{X}}$ to zero is a consequence of Proposition 36 and of the Borel–Cantelli lemma. It occurs under the assumption $\rho > 1$. ■

C.6 The power of the test

Proof of Proposition 21

Lemma 39.

Let α, κ be two positive numbers and \mathcal{L} and \mathcal{L}^* two laws of real random variables. We denote q_α (respectively q_α^*) the α -quantile of the law \mathcal{L} (respectively \mathcal{L}^*). If $W_1(\mathcal{L}, \mathcal{L}^*) < \kappa$ then:

$$q_\alpha^* \leq 2\frac{\kappa}{\alpha} + q_{\frac{\alpha}{2}}.$$

Proof

With a drawing, since the L_1 -norm between $F_{\mathcal{L}}$ and $F_{\mathcal{L}^*}$ is smaller than κ , we have:

$$F_{\mathcal{L}^*} \left(q_{\frac{\alpha}{2}} + 2\frac{\kappa}{\alpha} \right) > 1 - \alpha.$$

■

In this part we assume that m is fixed in $[0, 1]$ and $N = cn^\rho$ for some $\rho > 1$ and $c > 0$.

Recall that our aim is to upper bound the type II error, that is:

$$\mathbb{P}_{(\mu,\nu)} \left(\sqrt{n}W_1(d_{\hat{\mu}_N,m}(\hat{\mu}_n), d_{\hat{\nu}_N,m}(\hat{\nu}_n)) < \hat{q}_\alpha \right).$$

For some $\kappa = n^\gamma$ with γ in $[0, \frac{1}{2})$ to be chosen later, we first upper bound the quantile \hat{q}_α with high probability.

As noticed in the proof of Lemma 16, the law of $\sqrt{n}W_1(d_{\mu,m}(\hat{\mu}_n), d_{\mu,m}(\hat{\mu}'_n))$ converges to $\mathcal{L}(\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1)$, there is also the convergence of the first moments. So, for n big enough, we have:

$$W_1(\mathcal{L}(\sqrt{n}W_1(d_{\mu,m}(\hat{\mu}_n), d_{\mu,m}(\hat{\mu}'_n))), \mathcal{L}(\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1)) \leq 1.$$

Then, under the assumption

$$W_1(\mathcal{L}(\sqrt{n}W_1(d_{\mu,m}(\hat{\mu}_n), d_{\mu,m}(\hat{\mu}'_n)), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N))) \leq \kappa,$$

we have

$$W_1(\mathcal{L}(\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N)) \leq \kappa + 1.$$

We can do the same thing for ν . Thus we get that for n big enough and under the previous assumptions:

$$W_1\left(\frac{1}{2}\mathcal{L}(\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1) + \frac{1}{2}\mathcal{L}(\|\mathbb{G}_{\nu,m} - \mathbb{G}'_{\nu,m}\|_1), \frac{1}{2}\mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N) + \frac{1}{2}\mathcal{L}_{N,n,m}^*(\hat{\nu}_N, \hat{\nu}_N)\right) \leq \kappa + 1.$$

And thanks to Lemma 39,

$$\hat{q}_\alpha \leq \tilde{q}_{\frac{\alpha}{2}} + 2\frac{\kappa + 1}{\alpha},$$

with \tilde{q}_α the α -quantile of the law $\frac{1}{2}\mathcal{L}(\|\mathbb{G}_{\mu,m} - \mathbb{G}'_{\mu,m}\|_1) + \frac{1}{2}\mathcal{L}(\|\mathbb{G}_{\nu,m} - \mathbb{G}'_{\nu,m}\|_1)$.

We need to remark that with similar arguments as for Lemma 28, we have:

$$\begin{aligned} & W_1(\mathcal{L}(\sqrt{n}W_1(d_{\mu,m}(\hat{\mu}_n), d_{\mu,m}(\hat{\mu}'_n)), \mathcal{L}_{N,n,m}^*(\hat{\mu}_N, \hat{\mu}_N))) \leq \\ & 2\sqrt{n}\mathcal{D}_\mu \|F_{d_{\mu,m}(\mu)} - F_{d_{\mu,m}(\hat{\mu}_N)}\|_{\infty, (0, \mathcal{D}_\mu)} + 2\frac{\sqrt{n}}{m}W_1(\mu, \hat{\mu}_N). \end{aligned}$$

Now remark that

$$\begin{aligned} & \sqrt{n}W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\hat{\nu}_N, m}(\hat{\nu}_n)) \geq \sqrt{n}W_1(d_{\mu, m}(\mu), d_{\nu, m}(\nu)) \\ & - \sqrt{n}W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\mu, m}(\mu)) - \sqrt{n}W_1(d_{\hat{\nu}_N, m}(\hat{\nu}_n), d_{\nu, m}(\nu)), \end{aligned}$$

but as well, thanks to Lemma 32, the definition of the L_1 -Wasserstein distance as the L_1 -norm between the cumulative distribution functions and to Proposition 5:

$$\begin{aligned} & \sqrt{n}W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\mu, m}(\mu)) \leq \\ & \frac{\sqrt{n}}{m}W_1(\mu, \hat{\mu}_N) + \sqrt{n}\mathcal{D}_{\mu, m} \|F_{d_{\mu, m}(\hat{\mu}_n)} - F_{d_{\mu, m}(\mu)}\|_{\infty, (0, \mathcal{D}_\mu)}, \end{aligned}$$

with $\mathcal{D}_{\mu, m}$ the diameter of the support of the measure $d_{\mu, m}(\mu)$. So, we can finally upper bound $\mathbb{P}_{(\mu, \nu)}(\sqrt{n}W_1(d_{\hat{\mu}_N, m}(\hat{\mu}_n), d_{\hat{\nu}_N, m}(\hat{\nu}_n)) < \hat{q}_\alpha)$ by

$$\begin{aligned} & \mathbb{P}\left(\sqrt{n}\mathcal{D}_\mu \|F_{d_{\mu, m}(\mu)} - F_{d_{\mu, m}(\hat{\mu}_N)}\|_{\infty, (0, \mathcal{D}_\mu)} \geq \frac{\kappa}{4}\right) + \\ & \mathbb{P}\left(\sqrt{n}\mathcal{D}_\nu \|F_{d_{\nu, m}(\nu)} - F_{d_{\nu, m}(\hat{\nu}_N)}\|_{\infty, (0, \mathcal{D}_\nu)} \geq \frac{\kappa}{4}\right) + \\ & 2\mathbb{P}\left(\frac{\sqrt{n}}{m}W_1(\mu, \hat{\mu}_N) \geq \frac{\kappa}{4}\right) + 2\mathbb{P}\left(\frac{\sqrt{n}}{m}W_1(\nu, \hat{\nu}_N) \geq \frac{\kappa}{4}\right) + \\ & \mathbb{P}\left(\|F_{d_{\mu, m}(\hat{\mu}_n)} - F_{d_{\mu, m}(\mu)}\|_{\infty, (0, \mathcal{D}_\mu)} \geq \frac{W_1(d_{\mu, m}(\mu), d_{\nu, m}(\nu))}{2\mathcal{D}_{\mu, m}} - \frac{\tilde{q}_{\frac{\alpha}{2}}}{2\mathcal{D}_{\mu, m}\sqrt{n}} - \frac{(4 + \alpha)\kappa + 4}{4\mathcal{D}_{\mu, m}\alpha\sqrt{n}}\right) + \\ & \mathbb{P}\left(\|F_{d_{\nu, m}(\hat{\nu}_n)} - F_{d_{\nu, m}(\nu)}\|_{\infty, (0, \mathcal{D}_\nu)} \geq \frac{W_1(d_{\mu, m}(\mu), d_{\nu, m}(\nu))}{2\mathcal{D}_{\nu, m}} - \frac{\tilde{q}_{\frac{\alpha}{2}}}{2\mathcal{D}_{\nu, m}\sqrt{n}} - \frac{(4 + \alpha)\kappa + 4}{4\mathcal{D}_{\nu, m}\alpha\sqrt{n}}\right). \end{aligned}$$

For all positive ϵ , for n big enough, remark that the sum of the last two terms can be bounded thanks to the DKW-Massart inequality [22], by

$$4 \exp\left(-\frac{W_1^2(d_{\mu, m}(\mu), d_{\nu, m}(\nu))}{(2 + \epsilon) \max\{\mathcal{D}_\mu^2, \mathcal{D}_\nu^2\}}n\right).$$

Remark also that thanks to the DKW-Massart inequality, the first term can be upper bounded by

$$2 \exp\left(-\frac{1}{8\mathcal{D}_\mu^2} cn^{\rho-1+2\gamma}\right).$$

The second term is similar. Thanks to Theorem 2 in [17], the third term is upper bounded by

$$c_1 \exp\left(-c_2 m^d n^{\rho+d\gamma-\frac{d}{2}}\right),$$

for some fixed constants c_1 and c_2 . The remaining terms are similar.

Since $\rho > 1$, we can choose a positive γ satisfying: $\gamma < \frac{1}{2}$, $\rho + d\gamma - \frac{d}{2} > 1$ and $\rho - 1 + 2\gamma > 1$. So the two last expressions are negligible in comparison to the first one.

So, for n big enough, $\mathbb{P}_{(\mu,\nu)}(\sqrt{n}W_1(d_{\hat{\mu}_N,m}(\hat{\mu}_n), d_{\hat{\nu}_N,m}(\hat{\nu}_n)) < \hat{q}_\alpha)$ is upper bounded by

$$4 \exp\left(-\frac{W_1^2(d_{\mu,m}(\mu), d_{\nu,m}(\nu))}{3 \max\{\mathcal{D}_{\mu,m}^2, \mathcal{D}_{\nu,m}^2\}} n\right).$$

■

C.7 Numerical illustrations

In this section, we give details on the simulations presented in Section 5. Recall that we consider the measure μ_ν , that is, the distribution of the random vector $(R \sin(vR) + 0.03N, R \cos(vR) + 0.03N')$ with R , N and N' independent random variables; N and N' from the standard normal distribution and R uniform on $(0, 1)$.

From the measure μ_{10} we get a N -sample $P = \{X_1, X_2, \dots, X_N\}$, where $N = 2000$. As well, we get a N -sample $Q = \{Y_1, Y_2, \dots, Y_N\}$ from the measure μ_{20} . It leads to the empirical measures $\hat{\mu}_{10,N}$ and $\hat{\mu}_{20,N}$. On Figure 3, we plot the cumulative distribution function of the measure $d_{\hat{\mu}_{10,N},m}(\hat{\mu}_{10,N})$, that is, the function F defined for all t in \mathbb{R} by the proportion of the X_i in P satisfying $d_{\hat{\mu}_{10,N},m}(X_i) \leq t$. It approximates the true cumulative distribution function associated to the DTM-signature $d_{\mu,m}(\mu)$. As well, we plot the cumulative distribution function of the measure $d_{\hat{\mu}_{20,N},m}(\hat{\mu}_{20,N})$. Observe that the signatures are different. Thus, for the choice of parameter $m = 0.05$, the DTM-signature discriminates well between the measures μ_{10} and μ_{20} .

In Figure 4, for $m = 0.05$ and $n = 20$, we first generate $N_{MC} = 1000$ independent realisations of the random variable $\sqrt{n}W_1(d_{\hat{\mu}_{10,N-n},m}(\hat{\mu}_{10,n}), d_{\hat{\mu}'_{10,N-n},m}(\hat{\mu}'_{20,n}))$, where $\hat{\mu}_{10,N}$ and $\hat{\mu}'_{10,N}$ are independent empirical measures from μ_{10} , $\hat{\mu}_{10,N} = \frac{n}{N}\hat{\mu}_{10,n} + \frac{N-n}{N}\hat{\mu}_{10,N-n}$ and $\hat{\mu}'_{10,N} = \frac{n}{N}\hat{\mu}'_{10,n} + \frac{N-n}{N}\hat{\mu}'_{10,N-n}$. We plot the empirical cumulative distribution function associated to this N -sample. As well, from two fixed N -samples from the law μ_{10} , P and Q , we generate a set *boot* of N_{MC} random variables, as explained in the Algorithm in Section 4.1, and we plot its cumulative distribution function. Remark that the two cumulative distribution functions are close. It means that the α -quantile of the distribution of the test statistic is well approximated by the α -quantile of the bootstrap distribution.

The Figure 5 is obtained by applying the test **DTM** and the test **KS** to two independent N -samples, 1000 times independently, and by averaging the number of rejections of the hypothesis H_0 . For the type-I error, the N -samples are both from μ_{10} , as for the power, a sample is from μ_{10} and the other one from μ_ν .