



# Recovering the homology of immersed manifolds

Raphaël Tinarrage

## ► To cite this version:

Raphaël Tinarrage. Recovering the homology of immersed manifolds. Discrete and Computational Geometry, 2023, 69 (3), pp.659-744. 10.1007/s00454-022-00409-5 . hal-02396261v4

**HAL Id: hal-02396261**

**<https://hal.science/hal-02396261v4>**

Submitted on 2 Apr 2024

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

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



Distributed under a Creative Commons Attribution 4.0 International License

# RECOVERING THE HOMOLOGY OF IMMERSED MANIFOLDS

Raphaël TINARRAGE

Datashape, Inria Paris-Saclay – LMO, Université Paris-Saclay

**Abstract.** Given a sample of an abstract manifold immersed in some Euclidean space, we describe a way to recover the singular homology of the original manifold. It consists in estimating its tangent bundle—seen as subset of another Euclidean space—from a measure theoretical point of view, and in applying measure-based filtrations for persistent homology. We show that our construction is consistent and stable. The proof relies on two main ingredients. First, we introduce and study the normal reach, a notion of reach adapted to immersed manifolds. It allows to quantify the deviation of geodesics around self-intersections. Secondly, we study the estimation of tangent spaces via local principal component analysis, with respect to the Wasserstein distance. We illustrate our method on a few synthetic datasets, in the context of homology estimation and transverse manifolds clustering.

**Numerical experiments.** A Python notebook can be found at <https://github.com/raphaeltinarrage/ImmersedManifolds/blob/master/Demo.ipynb>. Some animations are gathered at <https://youtube.com/playlist?list=PLFkltNTtklDlIFg1djM5XprlL8Ys0hW4>.

**MSC codes.** 55N31, 53C42, 53C20, 49Q15, 49Q22, 68U05.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Preliminaries</b>	<b>5</b>
2.1	Euclidean and Riemannian geometry . . . . .	5
2.2	Persistent homology . . . . .	9
2.3	Persistent homology for measures . . . . .	10
2.4	Model and hypotheses . . . . .	11
<b>3</b>	<b>Reach of an immersed manifold</b>	<b>13</b>
3.1	Normal reach . . . . .	14
3.2	Probabilistic bounds under normal reach conditions . . . . .	17
3.3	Sublevel sets of the normal reach . . . . .	24
<b>4</b>	<b>Tangent space estimation</b>	<b>30</b>
4.1	Local covariance matrices and lifted measures . . . . .	30
4.2	Consistency of the estimation . . . . .	32
4.3	Stability of localization of measures . . . . .	35
4.4	Stability of the estimation . . . . .	45
4.5	An approximation theorem . . . . .	51
<b>5</b>	<b>Topological inference with the lifted measure</b>	<b>52</b>
5.1	Overview of the method . . . . .	52
5.2	Homotopy type estimation with the DTM . . . . .	55
5.3	Persistent homology with DTM-filtrations . . . . .	57
<b>6</b>	<b>Conclusion</b>	<b>61</b>
<b>A</b>	<b>Notations</b>	<b>61</b>
<b>B</b>	<b>Table of constants</b>	<b>62</b>
<b>C</b>	<b>Supplementary material for Sect. 2</b>	<b>63</b>

# 1 Introduction

A central challenge in Topological Data Analysis (TDA) consists in estimating the topology of a subset  $\mathcal{M} \subset \mathbb{R}^n$  based on a finite collection of points  $X$  that lie in  $\mathcal{M}$  or close to. By estimating the topology of  $\mathcal{M}$ , we mean inferring its homotopy type, or more simply inferring its singular homology groups. In what follows, the subset  $\mathcal{M}$  will be referred to as the *underlying space*, and  $X$  as the *observation*.

Inferring the homotopy type of  $\mathcal{M}$  may be done by constructing a homotopy equivalent simplicial complex. A usual method consists in considering the union of balls of radius  $t \geq 0$  centered around every point of  $X$ , and in taking the nerve of this covering [1]. This simplicial complex is called the *Čech complex* of  $X$  with parameter  $t$ . One can also consider the *Vietoris-Rips complex* of  $X$  with parameter  $t$ , defined as the clique complex of the underlying graph of the previous complex. The parameter  $t$  is to be chosen in accordance with the Hausdorff distance  $d_H(X, \mathcal{M})$  and some geometric quantities associated to  $\mathcal{M}$ , such as its reach [2, 3, 4] or its  $\mu$ -reach [5, 6]. Several variations of this construction have been studied, for instance by letting the parameter  $r$  vary across the points of  $X$  [3, 4], by considering ellipsoids instead of balls [7], or by using balls restricted to  $\mathcal{M}$  [4]. Besides the Čech and the Rips complex, one may also consider the  $\alpha$ -shape, obtained by first building the Delaunay triangulation of  $X$ , and then keeping simplices that fit in an empty ball of radius  $\alpha$ . This construction yields a simplicial complex homotopy equivalent to the Čech complex [8, 9]. Developments of this construction include the *witness complex* [10, 11], obtained by choosing a subset of ‘landmark’ points, or the *tangential Delaunay complex* [12], that incorporates tangent space information.

Besides, the problem of inference of homology groups of  $\mathcal{M}$  can be solved by computing a homotopy equivalent simplicial complex, such as those listed in the previous paragraph. However, other solutions to this problem have been proposed. They often consist in computing the image of the map induced in simplicial homology by a simplicial inclusion  $K^s \hookrightarrow K^t$ , where  $K^s$  (resp.  $K^t$ ) is the Čech or the Vietoris-Rips complex at time  $s$  (resp.  $t$ ). The parameters  $s$  and  $t$  are still to be chosen in accordance with the Hausdorff distance  $d_H(X, \mathcal{M})$  and some geometric quantities of  $\mathcal{M}$ , such as its weak feature size [13, 14] or its convexity radius and distortion [15].

Another point of view on inference of homology groups, that allows to avoid the selection of the parameters  $s$  and  $t$ , is *persistent homology* [16, 17]. It consists in building from  $X$  an algebraic structure, called a *persistence module*, which can be summarized in a *persistence barcode*. The bars of the barcode can be interpreted as homological features of  $X$  at different scales. These persistence modules are obtained from *filtrations*, that is, increasing families of subspaces built on top of  $X$ . Among the many filtrations available to the user, the most used are the sublevel sets of the distance function to  $X$ , its simplicial equivalent the *Čech filtration*, and its clique-complex version the *Vietoris-Rips filtration*. The main theoretical advantage of these filtrations is their stability: small perturbations of  $X$  in Hausdorff distance implies only small perturbations of the barcodes in bottleneck distance [18]. This stability allows to design statistical procedures for inferring the homology groups of  $\mathcal{M}$  from  $X$  [19, 20, 21].

A critical problem, both in the context of homotopy type inference and homology inference, is the presence of *anomalous points* in  $X$ , that is, roughly speaking, points that cause the Hausdorff distance  $d_H(X, \mathcal{M})$  to be large. In presence of anomalous points, the results presented above cannot be used. Among the attempts that have been made to overcome this issue, the filtration defined by the sublevel sets of the *distance-to-measure* (DTM) introduced in [22], and some of its variants [23], have been proven to provide relevant information. Unfortunately, from a practical perspective, the exact computation of the sublevel sets filtration of the DTM turn out to be far too expensive in most cases. To address this problem, the *witnessed  $k$ -distance* [24], the *weighted Vietoris-Rips complex filtration* [25] and the *DTM-filtrations* [26] have been proposed.

In this paper, we address the problem of homotopy type and homological inference, by weakening the assumptions of [2], where it is supposed that  $\mathcal{M}$  is a submanifold with positive reach. Here, we consider that  $\mathcal{M}$  is an immersed manifold, not embedded. That is to say, we suppose that there exists an abstract  $\mathcal{C}^2$ -manifold  $\mathcal{M}_0$ , immersed in the Euclidean space via a  $\mathcal{C}^2$ -immersion  $u: \mathcal{M}_0 \rightarrow \mathbb{R}^n$ , whose image is  $\mathcal{M}$ . As before, the observation  $X$  is a subset of  $\mathbb{R}^n$ , that we suppose close to  $\mathcal{M}$  in Hausdorff distance. Throughout this paper, we will use the example of a circle immersed in the plane as a lemniscate, as represented in Figure 1. Being an immersion,  $\mathcal{M}$  may self-intersect, and the sets  $\mathcal{M}_0$  and  $\mathcal{M}$  may have different homotopy types. The Čech filtration of  $\mathcal{M}$ , or  $X$ , would reveal the homology of  $\mathcal{M}$ , not that

of  $\mathcal{M}_0$ . Consequently, the usual approach based on the Čech filtration no longer applies here, and new methods must be developed.

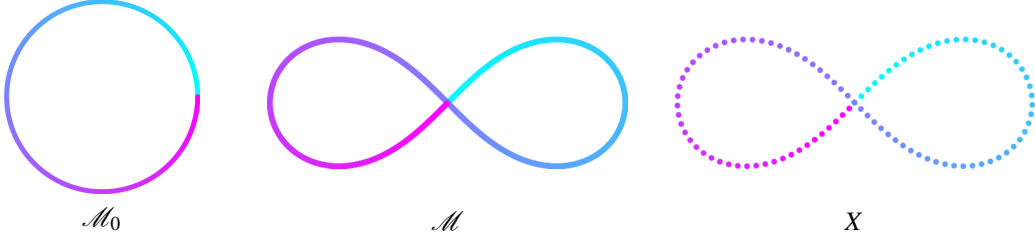


Figure 1: Left: The abstract manifold  $\mathcal{M}_0$ , a circle. Middle: The immersion  $\mathcal{M} \subset \mathbb{R}^2$ , known as the lemniscate of Bernoulli. Right: The observation  $X$ .

**Previous work.** Among the works that involve immersed manifolds, let us cite [27, 28, 29], which are set in the context where  $\mathcal{M}$  is a union of intersecting submanifolds. Hence  $\mathcal{M}$  is not a submanifold itself, but it is an immersed manifold, coming from an abstract manifold  $\mathcal{M}_0$ , made up of several connected components. In these three works, the authors propose algorithms to classify the different components of  $\mathcal{M}$ . In the context of the present paper, classifying the components of  $\mathcal{M}$  means finding the connected components of  $\mathcal{M}_0$ . Each of these algorithms rely on the estimation of tangent spaces, so as to separate the set  $\mathcal{M}$  where it self-intersects. In other words, they estimate the tangent bundle of the manifold. This is a point of view that we also adopt. We remark that, among these works, only [29] provides mathematical proofs of consistency, for their Algorithms 2 and 3. We compare this method to ours at the end of this subsection.

Another related problem is the one of dimension estimation. In many manifold reconstruction algorithms that involve the estimation of tangent spaces, such as in [27, 12, 30, 31, 32], or in [29, Algorithm 4], the dimension  $d$  of the underlying manifold  $\mathcal{M}$  is given as an input of the algorithm. If  $d$  is not known, a dimension estimator may be used, whether supposing that the input data exactly lies on  $\mathcal{M}$  [33, 34], or allowing the data to be corrupted by noise [35, 36, 37, 28]. Another strategy consists in designing tangent spaces estimators that does not require the dimension  $d$ , such as the empirical covariance matrix [29, Algorithms 2 and 3]. In the present paper, we generalize the definition of the empirical covariance matrix to any measure input, that we call *local covariance matrices* (see Definition 4.1). We show that it is a consistent estimator of the tangent spaces (Proposition 4.1) and that is is robust to noise (see Equation (40)).

Our method is based on the stability of tangent space estimation via local covariance matrices. Such a stability has already been studied in [38], and the stability of truncations of measures in [39].

**Our contributions.** In order to estimate the homology of a manifold from an immersion of it, we propose to estimate its tangent bundle, seen as a subset of another Euclidean space. As it turns out, in the process of estimating this tangent bundle, we will make errors, which will result in anomalous points. This issue will be solved by using the DTM-filtrations, which require to use a measure theoretical framework [22, 26]. Let us describe the method, in measure theoretical terms.

Let  $\mathcal{M}_0$  be a compact  $\mathcal{C}^2$ -manifold of dimension  $d$ , and  $\mu_0$  a Radon probability measure on  $\mathcal{M}_0$  with full support. Let  $u: \mathcal{M}_0 \rightarrow \mathbb{R}^n$  be a  $\mathcal{C}^2$ -immersion. We assume the following genericity condition: the immersion is such that self-intersection points correspond to different tangent spaces. In other words, for every  $x_0, y_0 \in \mathcal{M}_0$  such that  $x_0 \neq y_0$  and  $u(x_0) = u(y_0)$ , the tangent spaces  $d_{x_0}u(T_{x_0}\mathcal{M}_0)$  and  $d_{y_0}u(T_{y_0}\mathcal{M}_0)$  of  $\mathcal{M}_0$ , seen in  $\mathbb{R}^n$ , are different. As we will explain later, this condition ensures that the problem is well-posed (see Hypothesis 1). Now, define the image of the immersion  $\mathcal{M} = u(\mathcal{M}_0)$  and the pushforward measure  $\mu = u_*\mu_0$ . We consider the following problem: the input data is the measure  $\mu$ , or a close measure  $\nu$ . Our goal is to infer the singular homology of  $\mathcal{M}_0$  (with coefficients in  $\mathbb{Z}/2\mathbb{Z}$  for instance) from the data  $\nu$ . In practice,  $\nu$  can be given as the empirical measure on a point cloud. To answer this problem, we will build in this paper a persistence module such that the homology of  $\mathcal{M}_0$  can be read on the corresponding persistence diagram.

To get back to  $\mathcal{M}_0$ , we proceed as follows: let  $M(\mathbb{R}^n)$  be the vector space of  $n \times n$  matrices, and  $\check{u}: \mathcal{M}_0 \rightarrow \mathbb{R}^n \times M(\mathbb{R}^n)$  the map

$$\check{u}: x_0 \mapsto \left( u(x_0), \frac{1}{d+2} p_{T_{u(x_0)}\mathcal{M}} \right),$$

where  $p_{T_{u(x_0)}\mathcal{M}}$  is the matrix of the orthogonal projection on the tangent space  $T_{u(x_0)}\mathcal{M} = d_{x_0}u(T_{x_0}\mathcal{M}_0) \subset \mathbb{R}^n$ , written in the canonical basis of  $\mathbb{R}^n$ . The term  $\frac{1}{d+2}$  is a technical normalization factor that will be explained later (see Proposition 4.1). Now, define the set  $\check{\mathcal{M}} = \check{u}(\mathcal{M}_0)$ . It is a submanifold of  $\mathbb{R}^n \times M(\mathbb{R}^n)$ ,  $\mathcal{C}^1$ -diffeomorphic to  $\mathcal{M}_0$ . It is called the *lift* of  $\mathcal{M}_0$ , or the *lifted manifold*. The space  $\mathbb{R}^n \times M(\mathbb{R}^n)$  is called the *lift space*. Figure 2 provides a representation of the lifted manifold, when the input immersion is the lemniscate, as in Figure 1.

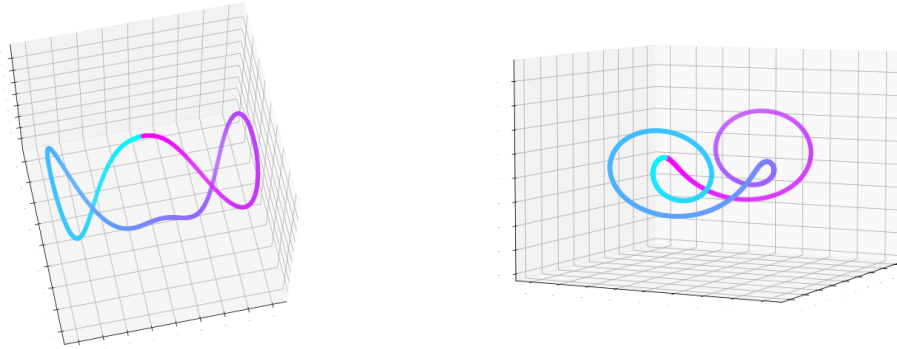


Figure 2: Two views of the submanifold  $\check{\mathcal{M}} \subset \mathbb{R}^2 \times M(\mathbb{R}^2) \simeq \mathbb{R}^6$ , projected in a 3-dimensional subspace via Principal Component Analysis (PCA). Observe that it does not self-intersect. The initial set  $\mathcal{M}$  is represented in Figure 1.

Suppose that one is able to estimate  $\check{\mathcal{M}}$  from  $v$ . Then one could consider the persistent homology of a filtration based on  $\check{\mathcal{M}}$ —say the Čech filtration of  $\check{\mathcal{M}}$  in the ambient space  $\mathbb{R}^n \times M(\mathbb{R}^n)$  for instance—and read the singular homology of  $\mathcal{M}_0$  in the corresponding persistent barcode. This is represented in Figure 3.



Figure 3: Left: Persistence barcode of the 1-homology of the Čech filtration of  $\mathcal{M}$  in the ambient space  $\mathbb{R}^2$ . One reads the 1-homology of the lemniscate. Right: Persistence barcode of the 1-homology of the Čech filtration of  $\check{\mathcal{M}}$  in the lift space  $\mathbb{R}^2 \times M(\mathbb{R}^2)$ . At the beginning of the barcode, one reads the 1-homology of a circle. Parameter  $\gamma = 2$ .

Unfortunately, we won't be able to give a good estimation of  $\check{\mathcal{M}}$ . This is because the tangent spaces  $T_{u(x_0)}\mathcal{M}$ , that we compute via local covariance matrices, won't be estimated correctly if  $x$  is too close to a self-intersection of  $\mathcal{M}$ . In order to get around this issue, we adopt a measure theoretical point of view. Instead of estimating the lifted submanifold  $\check{\mathcal{M}}$ , we propose to estimate the *exact lifted measure*  $\check{\mu}_0$ , defined as the push-forward  $\check{\mu}_0 = \check{u}_*\mu_0$ . It is a measure on the lift space  $\mathbb{R}^n \times M(\mathbb{R}^n)$  and has support  $\check{\mathcal{M}}$ .

It is worth noting that  $\check{\mathcal{M}}$  can be naturally seen as a submanifold of  $\mathbb{R}^n \times \mathcal{G}_d(\mathbb{R}^n)$ , where  $\mathcal{G}_d(\mathbb{R}^n)$  denotes the Grassmannian of  $d$ -dimensional linear subspaces of  $\mathbb{R}^n$ . From this point of view,  $\check{\mu}_0$  can be seen as a measure on  $\mathbb{R}^n \times \mathcal{G}_d(\mathbb{R}^n)$ , i.e., a *varifold*. This point of view has already been used in

data analysis, such as in geometric inference [40, 41] or in computational anatomy [42]. However, for computational reasons, we choose to work in the matrix space  $\mathbf{M}(\mathbb{R}^n)$  instead of  $\mathcal{G}_d(\mathbb{R}^n)$ .

Here is an alternative definition of  $\check{\mu}_0$ : for any test function  $\phi : \mathbb{R}^n \times \mathbf{M}(\mathbb{R}^n) \rightarrow \mathbb{R}$ ,

$$\int \phi(x, A) d\check{\mu}_0(x, A) = \int_{\mathcal{M}_0} \phi \left( u(x_0), \frac{1}{d+2} p_{T_{u(x_0)} \mathcal{M}} \right) d\mu_0(x_0).$$

Getting back to the observed measure  $\nu$ , we propose to estimate  $\check{\mu}_0$  with the *lifted measure*  $\check{\nu}$ , defined as follows: for any test function  $\phi : \mathbb{R}^n \times \mathbf{M}(E) \rightarrow \mathbb{R}$ ,

$$\int \phi(x, A) d\check{\nu}(x, A) = \int_{\mathcal{M}} \phi \left( x, \bar{\Sigma}_\nu(x) \right) d\nu(x),$$

where  $\bar{\Sigma}_\nu(x)$  is *normalized local covariance matrix* (see Definition 4.1). It depends on a parameter  $r > 0$ . We prove that  $\bar{\Sigma}_\nu(x)$  can be used to estimate the tangent spaces  $\frac{1}{d+2} p_{T_{u(x_0)} \mathcal{M}}$  of  $\mathcal{M}$ . However, this estimation is biased next to the self-intersection of  $\mathcal{M}$ , as shown in Figure 4. As a consequence, the support of  $\check{\nu}$  is not close to  $\mathcal{M}$  in Hausdorff distance.

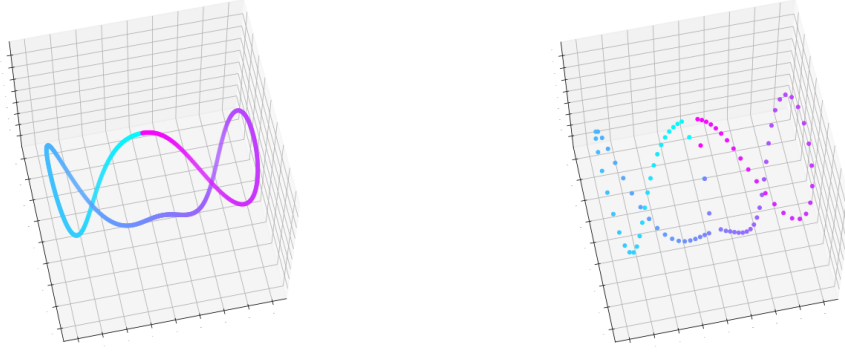


Figure 4: Left: The set  $\text{supp}(\check{\mu}_0) = \mathcal{M}$ , where  $\mu$  is the uniform measure on  $\mathcal{M}$  (see Figure 1). Right: The set  $\text{supp}(\check{\nu})$ , where  $\nu$  is the empirical measure on  $X$ . Parameters  $\gamma = 2$  and  $r = 0.1$ .

At this point, one could use an outliers-removal procedure, so as to recover  $\mathcal{M}$ . However, such a procedure depends critically on a choice of parameter, and is not reliable in practice. Instead, and still from a measure theoretical point of view, we will prove that the measure  $\check{\nu}$  is close to  $\check{\mu}_0$  in *Wasserstein distance* (see Theorem 4.14). This is true since only a few anomalous points are present. As a consequence, by using persistent homology for measures—such as the DTM-filtrations—the measure  $\check{\nu}$  can be used to infer the homotopy type of  $\mathcal{M}$ , that is, of  $\mathcal{M}_0$  (see Corollaries 5.3 and 5.5). The barcodes of the DTM-filtration on  $\check{\nu}$  are represented in Figure 5.

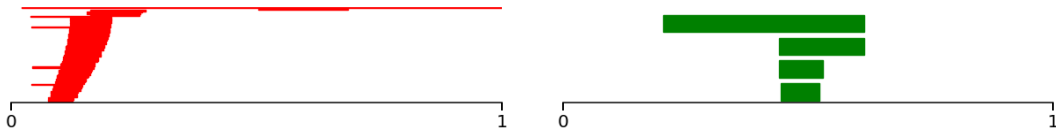


Figure 5: Persistence barcodes of the 0-homology (left) and 1-homology (right) of the DTM-filtration of the lifted measure  $\check{\nu}$ . Observe that the homology of the circle is salient on these barcodes (one large red bar and one large green bar). Parameters  $\gamma = 2$ ,  $r = 0.1$  and  $m = 0.01$ .

In order to quantify the quality of this approximation, we introduce a new geometric quantity: the *normal reach* (see Definition 3.1). It has been designed to play the role of the reach, when the subset considered is an immersed manifold. We show that the normal reach gives a scale at which an immersed manifold can be seen as an embedded manifold (see Proposition 3.8).

As a last remark, let us compare our method to [29]. In this paper, the input dataset is a point cloud  $X \subset \mathbb{R}^n$ , seen as a sample of the union of two intersecting submanifolds. Translated in our context,  $\mathcal{M}_0$  is the disjoint union of two abstract manifolds, and  $\mathcal{M}$  is an immersion of it. Their Algorithm 3 consists in estimating the tangent spaces  $Q_x$  on top of each point  $x \in X$ , via a variation of the empirical covariance matrix. Then, the authors build a graph  $G$ , whose vertices are the input data points  $x \in X$ , and where an edge  $[x, y]$  is added if the Euclidean positions are close enough ( $\|x - y\| \leq \varepsilon$ ) and if the tangent space estimations are close enough too ( $\|Q_x - Q_y\|_F \leq \eta$ ). The output of the algorithm is then the connected components of  $G$ . Unfortunately, due to the bad estimation of tangent spaces around self-intersections, the algorithm may treat the intersection points as a cluster of its own, hence returning more connecting components than wanted. In order to circumvent this issue, their Algorithm 2 includes an outliers-removal step, so as to exclude points close to the self-intersection. Under a particular choice of parameters, it is shown that the algorithm returns exactly two clusters, accurately clustering the points away from the intersection.

In comparison, our method has been thought to estimate the singular homology of  $\mathcal{M}_0$ , not only its connected components. In this setting, the outliers-removal procedure is a crucial step. This is because removing too many points would cause the apparition of gaps in the lifted manifold  $\tilde{\mathcal{M}}$ , which would be complicated to fill. Instead of discarding outliers, our method incorporates a sort of hierarchical clustering, performed by the use of the DTM-filtrations. Indeed, in the DTM-filtration of the lifted measure  $\tilde{\mathbf{V}}$ , the points are weighted according to their degree of anomalousness. This anomalousness is quantified via their local density in the lift space  $\mathbb{R}^n \times \mathbf{M}(\mathbb{R}^n)$ . The underlying idea is the following: since only a few points are close to the intersection, only a few points will have a bad tangent space estimation, hence their density will be small. A careful analysis will make this idea rigorous. Another advantage of our method lies in the use of persistent homology: the output of our algorithm is a persistence barcode. Hence we do not need to select precise connected components, or more generally, precise homological features. It is up to the user to read on this barcode the bars that seem to be relevant (in general, one chooses the longest bars). This procedure is justified theoretically by Corollary 5.5, which shows that the output barcode is stable.

**Data availability.** A Python notebook, containing numerical illustrations and codes used in this paper, can be found at <https://raphaeltinarrage.github.io/ImmersedManifolds>.

**Outline.** The rest of the paper is as follows. Sect. 2 gathers usual definitions related to Euclidean topology of compact sets, Riemannian geometry and persistent homology. We also describe our model. In Sect. 3 we introduce the normal reach, and derive certain probability bounds based on it. In Sect. 4, we study the tangent space estimation of an immersed manifold via local covariance matrices. We gather these results in Sect. 5 to obtain estimation guarantees for our method.

**Notations and constants.** We gather in Appendix A the notations that are used. Moreover, throughout the paper, we will refer to constants that are collected in a table in Appendix B. It is not necessary to read this table, since the constants will be introduced along the text.

## 2 Preliminaries

### 2.1 Euclidean and Riemannian geometry

In this subsection, we give some geometry results that will be useful in what follows. Here and in the rest of the paper, we will only consider compact manifolds and submanifolds without boundary, and measures that are Radon measures. We refer the reader to [43] for an exposition of the notion of reach, to [44] for a presentation of Riemannian geometry, and to [45] for a gentle introduction to geometric measure theory.

**Reach.** Let  $X$  be any subset of  $\mathbb{R}^n$  and  $y \in \mathbb{R}^n$  a point. The *distance* from  $y$  to  $X$  is the quantity

$$\text{dist}(y, X) = \inf\{\|x - y\| \mid x \in X\}.$$

A *projection* of  $y$  on  $X$  is a point  $x \in X$  that minimizes the distance  $\|x - y\|$ . The *medial axis* of  $X$  is the subset  $\text{med}(X) \subset \mathbb{R}^n$  which consists of points  $y \in \mathbb{R}^n$  that admit at least two distinct projections on  $X$ :

$$\text{med}(X) = \{y \in \mathbb{R}^n \mid \exists x, x' \in X, x \neq x', \|y - x\| = \|y - x'\| = \text{dist}(y, X)\}.$$

The *reach* of  $X$  is

$$\text{reach}(X) = \inf \{\|x - y\| \mid x \in X, y \in \text{med}(X)\}.$$

A useful property of sets with positive reach is the approximation by tangent spaces. For a general set  $X$ , we define its tangent cone at  $x \in X$ , denoted  $\text{Tan}(X, x)$ , as:

$$\{0\} \cup \left\{ v \in \mathbb{R}^n \mid \forall \varepsilon > 0, \exists y \in X \text{ s.t. } y \neq x, \|y - x\| < \varepsilon, \left\| \frac{v}{\|v\|} - \frac{y - x}{\|y - x\|} \right\| < \varepsilon \right\}.$$

Note that if  $X$  is a submanifold, we recover the usual notion of tangent space. The following characterization is fundamental in the study of sets with positive reach:

**Theorem 2.1** ([43, Theorem 4.18(2)]). *A closed set  $X \subset \mathbb{R}^n$  has positive reach  $\tau$  if and only if for every  $x, y \in X$ , we have*

$$\text{dist}(y - x, \text{Tan}(X, x)) \leq \frac{1}{2\tau} \|y - x\|^2.$$

The reach is a quantity that controls both the local and global regularity of the set  $X$ . When  $X = \mathcal{M}$  is a topological submanifold, having a positive reach implies that  $\mathcal{M}$  is of regularity  $\mathcal{C}^{1,1}$  [46, Proposition 1.4]. Conversely, a  $\mathcal{C}^{1,1}$ -submanifold  $\mathcal{M}$  has a positive reach [43, Theorem 4.19]. Moreover, when  $\mathcal{M}$  is  $\mathcal{C}^2$ , it can be shown that  $\text{reach}(\mathcal{M})$  is caused either by a bottleneck structure or by high curvature:

**Theorem 2.2** ([31, Theorem 3.4]). *A closed  $\mathcal{C}^2$ -submanifold  $\mathcal{M} \subset \mathbb{R}^n$  with positive reach must satisfy at least one of the following two properties:*

- Global case: *there exist  $x, y \in \mathcal{M}$  with  $\|x - y\| = 2\text{reach}(\mathcal{M})$  and  $\frac{1}{2}(x + y) \in \text{med}(\mathcal{M})$ ,*
- Local case: *there exists an arc-length parametrized geodesic  $\gamma: I \rightarrow \mathcal{M}$  with  $\|\ddot{\gamma}(0)\| = \text{reach}(\mathcal{M})^{-1}$ .*

In this paper, we will suppose that the manifold is of regularity  $\mathcal{C}^2$ , so as to obtain uniform bounds on its second derivatives (see Hypothesis 2). We do not study whether the results could be generalized to  $\mathcal{C}^{1,1}$  manifolds.

**Riemannian structure on immersed manifolds.** If  $u: \mathcal{M}_0 \rightarrow \mathcal{M} \subset \mathbb{R}^n$  is an immersion of a  $\mathcal{C}^2$ -manifold, then  $\mathcal{M}_0$  is naturally endowed with a Riemannian structure, by pulling back the inner product of  $\mathbb{R}^n$ . This makes  $u$  an isometry. From now on, we will consider that  $\mathcal{M}_0$  is given this Riemannian structure. We denote the (abstract) tangent space of  $\mathcal{M}_0$  at  $x_0$  as  $T_{x_0}\mathcal{M}_0$ , its image in  $\mathbb{R}^n$  as  $T_{u(x_0)}\mathcal{M} = d_{x_0}u(T_{x_0}\mathcal{M}_0)$ , and its orthogonal complement, the normal space, as  $(T_{u(x_0)}\mathcal{M})^\perp$ . The geodesic distance between two points  $x_0, y_0 \in \mathcal{M}_0$  is denoted  $d_{\mathcal{M}_0}(x_0, y_0)$ . For any  $x_0 \in \mathcal{M}_0$  and  $r \geq 0$ , we denote by  $\mathcal{B}_{\mathcal{M}_0}(x_0, r)$  (resp.  $\overline{\mathcal{B}}_{\mathcal{M}_0}(x_0, r)$ ) the open (resp. closed) geodesic ball of center  $x_0$  and radius  $r$  of  $\mathcal{M}_0$ . Moreover, for any  $v_0 \in T_{x_0}\mathcal{M}_0$ , we denote by  $\mathcal{B}_{T_{x_0}\mathcal{M}_0}(v_0, r)$  the open ball of center  $v_0$  and radius  $r$  of  $T_{x_0}\mathcal{M}_0$ .

For every  $x_0 \in \mathcal{M}_0$ , one defines the *second fundamental form* of  $\mathcal{M}_0$  at  $x_0$ . It is a symmetric bilinear form

$$\Pi_{x_0}: T_{x_0}\mathcal{M}_0 \times T_{x_0}\mathcal{M}_0 \longrightarrow (T_{u(x_0)}\mathcal{M})^\perp.$$

Let  $x_0 \in \mathcal{M}_0$ ,  $v_0 \in T_{x_0}\mathcal{M}_0$  a unit vector, and consider an unit-speed geodesic  $\gamma_0: I \rightarrow \mathcal{M}_0$  such that  $\gamma_0(0) = x_0$  and  $\dot{\gamma}_0(0) = v_0$ . Let us denote by  $\gamma$  the map  $u \circ \gamma_0: I \rightarrow \mathcal{M}$ . The following relation can be found in [2, Sect. 6] or [47, Sect. 3]:

$$\Pi_{x_0}(v_0, v_0) = \ddot{\gamma}(0). \quad (1)$$



In particular, any bound on the operator norm  $\|\Pi_{x_0}\|_{\text{op}}$  of  $\Pi_{x_0}$  implies a bound on  $\|\dot{\gamma}(0)\|$ . From now, we suppose that the operator norms  $\|\Pi_{x_0}\|_{\text{op}}$  are bounded by a constant  $\rho > 0$  (see Hypothesis 1). For instance, if  $\mathcal{M}_0$  is an embedded manifold, then  $\rho$  can be chosen as its reach [2, Proposition 6.1]. In general, if  $\mathcal{M}_0$  is a compact  $\mathcal{C}^2$ -manifold, such a global upper bound  $\rho$  exists. Let us list a few useful results.

**Lemma 2.3.** *Let  $x_0 \in \mathcal{M}_0$  and  $\gamma_0: I \rightarrow \mathcal{M}_0$  an arc-length parametrized geodesic starting from  $x_0$ . Let  $\gamma = u \circ \gamma_0$ ,  $v = \dot{\gamma}(0)$  and  $x = u(x_0)$ . For all  $t \in I$ , we have*

1.  $\|\gamma(t) - (x + tv)\| \leq \frac{\rho}{2}t^2$ .

Consequently, for every  $y_0 \in \mathcal{M}_0$ , denoting  $\delta = d_{\mathcal{M}_0}(x_0, y_0)$  and  $y = u(y_0)$ , we have

2.  $\text{dist}(y - x, T_x\mathcal{M}) \leq \frac{\rho}{2}\delta^2$ ,
3.  $(1 - \frac{\rho}{2}\delta)\delta \leq \|x - y\|$ .

Concerning the immersion  $u: \mathcal{M}_0 \rightarrow \mathcal{M}$ , we deduce that

4. the map  $u$  is injective on the open geodesic ball  $\mathcal{B}_{\mathcal{M}_0}(x_0, \frac{2}{\rho})$ ,
5. for every  $y_0 \in \mathcal{B}_{\mathcal{M}_0}(x_0, \frac{1}{\rho})$  such that  $y_0 \neq x_0$ , the vector  $y - x$  is not orthogonal to  $T_x\mathcal{M}$  nor  $T_y\mathcal{M}$ .

The first point of this lemma can be found in [2, Equation (5)], and the other points follow directly. Note that stronger versions of these results can be found in [47].

We now state a technical lemma. It gives how much time it takes for a geodesic to exit a Euclidean ball (represented in Figure 6). It is a direct consequence of Lemma 2.3 and its proof is deferred to Appendix C.

**Lemma 2.4.** *Let  $x_0, y_0 \in \mathcal{M}_0$  and  $\gamma_0: I \rightarrow \mathcal{M}_0$  an arc-length parametrized geodesic with  $\gamma_0(0) = y_0$ . Define  $x = u(x_0)$ ,  $y = u(y_0)$ ,  $\gamma = u \circ \gamma_0$ ,  $v = \dot{\gamma}(0)$  and  $l = \|y - x\|$ . Suppose that  $l < \frac{1}{\rho}$  and  $\langle v, y - x \rangle = 0$ .*

1. The map  $t \mapsto \|\gamma(t) - x\|$  is increasing on  $[0, T_1]$  where  $T_1 = \frac{\sqrt{2}}{\rho} \sqrt{2 - \sqrt{3 + \rho^2 l^2}}$ .

Let  $r$  be such that  $l \leq r < \frac{1}{2\rho}$  and define

$$T_2 = \frac{\sqrt{2}}{\rho} \sqrt{1 - \rho r - \sqrt{1 - 2\rho r + \rho^2 l^2}} \quad \text{and} \quad T_2' = \frac{\sqrt{2}}{\rho} \sqrt{1 - \rho r + \sqrt{1 - 2\rho r + \rho^2 l^2}}.$$

2. If  $t \in (T_2, T_2')$ , then  $\|\gamma(t) - x\| > r$ . Moreover,  $T_2 \leq 2\sqrt{r^2 - l^2}$ .
3. If  $l = 0$ , then  $T_2 = \frac{1}{\rho}(1 - \sqrt{1 - 2\rho r})$  and  $T_2' = \frac{1}{\rho}(1 + \sqrt{1 - 2\rho r})$ .

Moreover, let  $s$  be such that  $0 \leq s \leq r$  and define

$$a = \inf\{t \geq 0 \mid \|\gamma(t) - x\| \geq s\} \quad \text{and} \quad b = \inf\{t \geq 0 \mid \|\gamma(t) - x\| \geq r\}.$$

4.  $b - a \leq \sqrt{6\sqrt{r^2 - s^2}}$ .
5. If  $l = 0$ , then  $b - a \leq 2(r - s)$ .

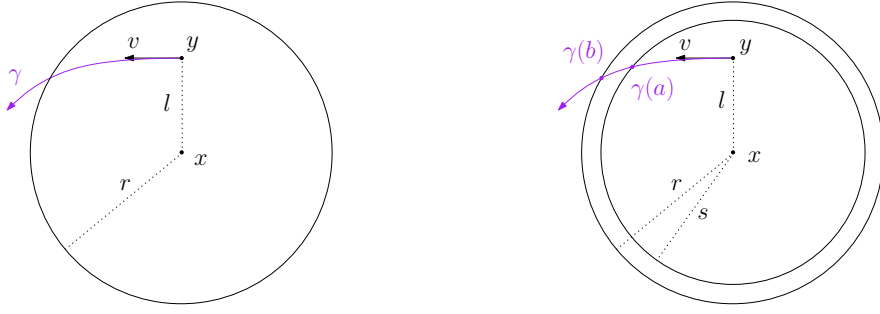


Figure 6: Illustration of Lemma 2.4 Point 1 (left) and Point 4 (right).

Last, the exponential map of  $\mathcal{M}_0$  at  $x_0$  will be denoted

$$\exp_{x_0}^{\mathcal{M}_0} : T_{x_0}\mathcal{M}_0 \rightarrow \mathcal{M}_0.$$

According to [48, Corollary 4 Point 1], the map  $\exp_{x_0}^{\mathcal{M}_0}$  is injective on the open ball  $\mathcal{B}_{T_{x_0}\mathcal{M}_0}\left(0, \frac{\pi}{\rho}\right)$  of  $T_{x_0}\mathcal{M}_0$ , and is a diffeomorphism onto its image  $\mathcal{B}_{\mathcal{M}_0}\left(0, \frac{\pi}{\rho}\right)$ . Moreover, for any  $x_0 \in \mathcal{M}_0$  and  $v_0 \in T_{x_0}\mathcal{M}_0$ , the  $d$ -dimensional Jacobian of  $\exp_{x_0}^{\mathcal{M}_0}$  at  $v_0$  is defined as

$$J_{v_0} = \sqrt{\det(A^t \cdot A)},$$

where  $A = d_{v_0} \exp_{x_0}^{\mathcal{M}_0}$  is the differential of the exponential map, seen as a  $d \times n$  matrix. As shown by the following result, the Jacobian of the exponential map is linked with the bound  $\rho$  on the operator norms  $\|\Pi_{x_0}\|_{\text{op}}$ .

**Lemma 2.5.** ([49, Proposition III.22]) *Let  $x_0 \in \mathcal{M}_0$  and  $v_0 \in T_{x_0}\mathcal{M}_0$  such that  $\|v_0\| = r < \frac{\pi}{2\sqrt{2}\rho}$ . The Jacobian  $J_{v_0}$  of  $\exp_{x_0}^{\mathcal{M}_0}$  at  $v_0$  satisfies*

$$\left(1 - \frac{(r\rho)^2}{6}\right)^d \leq J_{v_0} \leq \left(1 + (r\rho)^2\right)^d.$$

**Coarea formula.** For any measure  $\tau$  on a probability space  $\Omega$  and any measurable map  $h: \Omega \rightarrow \Omega'$ , the push-forward of  $\tau$  by  $h$  is the measure  $h_*\tau$  defined via

$$h_*\tau(A) = \tau(h^{-1}(A))$$

for any measurable set  $A \subset \Omega'$ . The transfer property refers to the following fact: for any integrable map  $\phi: \Omega' \rightarrow \mathbb{R}$ , we have

$$\int \phi \cdot dh_*\tau = \int \phi \circ h \cdot d\tau.$$

Now, suppose that  $\mathcal{M}_0$  and  $\mathcal{N}_0$  are Riemannian manifolds, of respective dimensions  $d \geq d'$ , and let  $\mathcal{H}_{\mathcal{M}_0}^d$  and  $\mathcal{H}_{\mathcal{N}_0}^{d'}$  denote their corresponding Hausdorff measures. The coarea formula allows to reformulate integrals on  $\mathcal{M}_0$  as integrals on  $\mathcal{N}_0$ .

**Theorem 2.6.** ([45, Chapter 3]) *Let  $f: \mathcal{M}_0 \rightarrow \mathcal{N}_0$  be a differentiable map. For any  $x_0 \in \mathcal{M}_0$ , let  $J_{x_0}$  denote its Jacobian. Let  $\phi: \mathcal{M}_0 \rightarrow [0, +\infty)$  be a measurable function. We have:*

$$\int_{\mathcal{M}_0} \phi(x_0) J_{x_0} \cdot d\mathcal{H}_{\mathcal{M}_0}^d(x_0) = \int_{\mathcal{N}_0} \left( \int_{x_0 \in f^{-1}(\{y_0\})} \phi(x_0) \cdot d\mathcal{H}_{\mathcal{M}_0}^{d-d'}(x_0) \right) d\mathcal{H}_{\mathcal{N}_0}^{d'}(y_0).$$

A useful consequence of this theorem is the following: suppose that  $\mathcal{M}_0$  is endowed with a Radon probability measure  $\mu_0$  that admits a density  $h: \mathcal{M}_0 \rightarrow [0, +\infty)$  against  $\mathcal{H}_{\mathcal{M}_0}^d$ . Consider the push-forward

measure  $\nu_0 = f_*\mu_0$ . If the Jacobian  $J_{x_0}$  of  $f$  never vanishes, then the push-forward measure  $\nu_0$  admits a density  $g: \mathcal{N}_0 \rightarrow [0, +\infty)$  against  $\mathcal{H}_{\mathcal{N}_0}^{d'}$ , where

$$g(y_0) = \int_{x_0 \in f^{-1}(\{y_0\})} h(x_0) J_{x_0}^{-1} \cdot d\mathcal{H}_{\mathcal{N}_0}^{d-d'}(x_0). \quad (2)$$

In particular, when  $d = d'$ , we have  $g(y_0) = \sum_{x_0 \in f^{-1}(\{y_0\})} h(x_0) J_{x_0}^{-1}$ .

## 2.2 Persistent homology

In this subsection, we write down the definitions of persistence modules, and their associated pseudo-distances, as presented in [18]. We refer the interested reader to [16, 17] for a thorough description. Let  $T \subset \mathbb{R}$  be an interval,  $E = \mathbb{R}^n$  a Euclidean space and  $k$  a field.

**Persistence modules.** A *persistence module* over  $T$  is a pair  $(\mathbb{V}, \mathbb{v})$  where  $\mathbb{V} = (V^t)_{t \in T}$  is a family of  $k$ -vector spaces, and  $\mathbb{v} = (v_s^t)_{s \leq t \in T}$  is a family of linear maps  $v_s^t: V^s \leftarrow V^t$  such that:

- for every  $t \in T$ ,  $v_t^t: V^t \rightarrow V^t$  is the identity map,
- for every  $r, s, t \in T$  such that  $r \leq s \leq t$ , we have  $v_s^t \circ v_r^s = v_r^t$ .

When there is no risk of confusion, we may denote a persistence module by  $\mathbb{V}$  instead of  $(\mathbb{V}, \mathbb{v})$ . Given  $\varepsilon \geq 0$ , an  $\varepsilon$ -*morphism* between two persistence modules  $\mathbb{V}$  and  $\mathbb{W}$  is a family of linear maps  $\phi = (\phi_t: V^t \rightarrow W^{t+\varepsilon})_{t \in T}$  such that the following diagram commutes for every  $s \leq t \in T$ :

$$\begin{array}{ccc} V^s & \xrightarrow{v_s^t} & V^t \\ \downarrow \phi_s & & \downarrow \phi_t \\ W^{s+\varepsilon} & \xrightarrow{w_{s+\varepsilon}^t} & W^{t+\varepsilon} \end{array}$$

If  $\varepsilon = 0$  and each  $\phi_t$  is an isomorphism, the family  $(\phi_t)_{t \in T}$  is an *isomorphism* of persistence modules. An  $\varepsilon$ -*interleaving* between two persistence modules  $\mathbb{V}$  and  $\mathbb{W}$  is a pair of  $\varepsilon$ -morphisms  $(\phi_t: V^t \rightarrow W^{t+\varepsilon})_{t \in T}$  and  $(\psi_t: W^t \rightarrow V^{t+\varepsilon})_{t \in T}$  such that the following diagrams commute for every  $t \in T$ :

$$\begin{array}{ccc} V^t & \xrightarrow{v_t^{t+2\varepsilon}} & V^{t+2\varepsilon} \\ & \searrow \phi_t & \nearrow \psi_{t+\varepsilon} \\ & W^{t+\varepsilon} & \end{array} \quad \begin{array}{ccc} & V^{t+\varepsilon} & \\ \psi_t \nearrow & & \searrow \phi_{t+\varepsilon} \\ W^t & \xrightarrow{w_t^{t+2\varepsilon}} & W^{t+2\varepsilon} \end{array}$$

The *interleaving pseudo-distance* between  $\mathbb{V}$  and  $\mathbb{W}$  is defined as

$$d_i(\mathbb{V}, \mathbb{W}) = \inf\{\varepsilon \geq 0 \mid \mathbb{V} \text{ and } \mathbb{W} \text{ are } \varepsilon\text{-interleaved}\}.$$

**Persistence barcodes.** A persistence module  $(\mathbb{V}, \mathbb{v})$  is said to be *pointwise finite-dimensional* if for every  $t \in T$ ,  $V^t$  is finite-dimensional. If this property is satisfied, we can define a notion of persistence barcode [50]. It comes from the algebraic decomposition of the persistence module into interval modules. Moreover, given two pointwise finite-dimensional persistence modules  $\mathbb{V}, \mathbb{W}$  with persistence barcodes  $\text{Barcode}(\mathbb{V})$  and  $\text{Barcode}(\mathbb{W})$ , the so-called isometry theorem states that

$$d_b(\text{Barcode}(\mathbb{V}), \text{Barcode}(\mathbb{W})) = d_i(\mathbb{V}, \mathbb{W}),$$

where  $d_i(\cdot, \cdot)$  denotes the interleaving distance between persistence modules, and  $d_b(\cdot, \cdot)$  denotes the bottleneck distance between barcodes.

More generally, the persistence module  $(\mathbb{V}, \mathbb{v})$  is said to be  $q$ -*tame* if for every  $s, t \in T$  such that  $s < t$ , the map  $v_s^t$  has finite rank. The  $q$ -tameness of a persistence module ensures that we can still define a notion of persistence barcode, even though the module may not be decomposable into interval modules. Moreover, the isometry theorem still holds [18].

**Filtrations of sets and simplicial complexes.** A family of subsets  $\mathbb{X} = (X^t)_{t \in T}$  of  $E$  is a *filtration* if it is non-decreasing for the inclusion, i.e., for any  $s, t \in T$  such that  $s \leq t$ , we have  $X^s \subset X^t$ . Given  $\varepsilon \geq 0$ , two filtrations  $\mathbb{X} = (X^t)_{t \in T}$  and  $\mathbb{Y} = (Y^t)_{t \in T}$  of  $E$  are  $\varepsilon$ -*interleaved* if, for every  $t \in T$ , we have  $X^t \subset Y^{t+\varepsilon}$  and  $Y^t \subset X^{t+\varepsilon}$ . The interleaving pseudo-distance between  $\mathbb{X}$  and  $\mathbb{Y}$  is defined as the infimum of such  $\varepsilon$ :

$$d_i(\mathbb{X}, \mathbb{Y}) = \inf\{\varepsilon \geq 0 \mid \mathbb{X} \text{ and } \mathbb{Y} \text{ are } \varepsilon\text{-interleaved}\}.$$

Filtrations of simplicial complexes and their interleaving distance are similarly defined: given a simplicial complex  $S$ , a *filtration of  $S$*  is a non-decreasing family  $\mathbb{S} = (S^t)_{t \in T}$  of subcomplexes of  $S$ . The interleaving pseudo-distance between two filtrations  $(S_1^t)_{t \in T}$  and  $(S_2^t)_{t \in T}$  of  $S$  is the infimum of the  $\varepsilon \geq 0$  such that they are  $\varepsilon$ -interleaved, i.e. for any  $t \in T$ , we have  $S_1^t \subset S_2^{t+\varepsilon}$  and  $S_2^t \subset S_1^{t+\varepsilon}$ .

**Relation between filtrations and persistence modules.** Applying the singular homology functor to a set filtration gives rise to a persistence module whose linear maps between homology groups are induced by the inclusion maps between sets. As a consequence, if two filtrations are  $\varepsilon$ -interleaved, then their associated persistence modules are also  $\varepsilon$ -interleaved, the interleaving homomorphisms being induced by the interleaving inclusion maps. As a consequence of the isometry theorem, if the modules are  $q$ -tame, then the bottleneck distance between their persistence barcodes is upper bounded by  $\varepsilon$  [18]. The same remarks hold when applying the simplicial homology functor to simplicial filtrations.

## 2.3 Persistent homology for measures

In this subsection we define the distance-to-measure (DTM), based on [22], and the DTM-filtrations, based on [26]. Let  $T = \mathbb{R}^+$  and  $E = \mathbb{R}^n$  endowed with the standard Euclidean norm.

**Wasserstein distances.** Given two probability measures  $\mu$  and  $\nu$  over  $E$ , a transport plan between  $\mu$  and  $\nu$  is a probability measure  $\pi$  over  $E \times E$  whose marginals are  $\mu$  and  $\nu$ . Let  $p \geq 1$ . The  $p$ -*Wasserstein distance between  $\mu$  and  $\nu$*  is defined as

$$W_p(\mu, \nu) = \left( \inf_{\pi} \int_{E \times E} \|x - y\|^p d\pi(x, y) \right)^{\frac{1}{p}},$$

where the infimum is taken over all the transport plans  $\pi$ . If  $q$  is such that  $p \leq q$ , then an application of Jensen's inequality shows that  $W_p(\mu, \nu) \leq W_q(\mu, \nu)$ .

**Distance-to-measure (DTM).** Let  $\mu$  be a probability measure over  $E$ , and  $m \in [0, 1)$  a parameter. The DTM associated to  $\mu$  with parameter  $m$  is the function  $d_{\mu, m}: E \rightarrow \mathbb{R}$  defined as:

$$d_{\mu, m}(x) = \sqrt{\frac{1}{m} \int_0^m \delta_{\mu, t}^2(x) dt} \quad \text{where} \quad \delta_{\mu, m}(x) = \inf\{r \geq 0 \mid \mu(\overline{\mathcal{B}}(x, r)) > m\},$$

and where  $\overline{\mathcal{B}}(x, r)$  denotes the closed ball of center  $x$  and radius  $r$  of  $E$ . When  $m$  is fixed and there is no risk of confusion, we may write  $d_\mu$  instead of  $d_{\mu, m}$ . Among the important properties of the DTM, it has been shown that it is 1-Lipschitz [22, Corollary 3.7]. Moreover, it is stable in Wasserstein distance [22, Theorem 3.5]: for any probability measures  $\mu$  and  $\nu$ , we have

$$\|d_{\mu, m} - d_{\nu, m}\|_\infty \leq m^{-\frac{1}{2}} W_2(\mu, \nu). \quad (3)$$

If  $f: E \rightarrow \mathbb{R}$  is any function and  $t \in \mathbb{R}$ , we will denote the  $t$ -*sublevel set* of  $f$  as  $f^t = f^{-1}((-\infty, t])$ . The following theorem shows that the sublevel sets  $d_{\mu, m}^t$  of  $d_{\mu, m}$  can be used to estimate the homotopy type of  $\text{supp}(\mu)$ .

**Theorem 2.7** ([22, Corollary 4.11, case  $\mu = 1$ ]). *Consider two probability measures  $\mu, \nu$  on  $E$  and  $m \in (0, 1)$ . Denote  $K = \text{supp}(\mu)$ . Suppose that  $\text{reach}(K) = \tau > 0$ , and that  $\mu$  satisfies the following hypothesis for  $r < (\frac{m}{a})^{\frac{1}{d}}$ :  $\forall x \in K, \mu(\overline{\mathcal{B}}(x, r)) \geq ar^d$ . Suppose that  $W_2(\mu, \nu) \leq m^{\frac{1}{2}}(\frac{\tau}{9} - (\frac{m}{a})^{\frac{1}{d}})$ . Define  $\varepsilon = (\frac{m}{a})^{\frac{1}{d}} + m^{-\frac{1}{2}} W_2(\mu, \nu)$  and choose  $t \in [4\varepsilon, \tau - 3\varepsilon]$ . Then  $d_{\nu, m}^t$  and  $K$  are homotopic equivalent.*

**DTM-filtrations.** These filtrations have been introduced in [26] and are defined as follows: consider a probability measure  $\mu$  on  $E$  and a parameter  $m \in [0, 1)$ . For every  $t \in T$ , consider the set

$$W^t[\mu] = \bigcup_{x \in \text{supp}(\mu)} \overline{\mathcal{B}}(x, t - d_{\mu, m}(x)), \quad (4)$$

where  $\overline{\mathcal{B}}(x, r)$  denotes the closed ball of center  $x$  and of radius  $r$  of  $E$  if  $r \geq 0$ , or denotes the empty set if  $r < 0$ . The family  $W[\mu] = (W^t[\mu])_{t \geq 0}$  is a filtration of  $E$ . It is called the *DTM-filtration* with parameters  $(\mu, m, 1)$ . By applying the singular homology functor, we obtain a persistence module, denoted  $\mathbb{W}[\mu]$ . If  $\text{supp}(\mu)$  is bounded, then  $\mathbb{W}[\mu]$  is  $q$ -tame. Moreover, it has been proven that the DTM-filtrations are stable with respect to the input measure:

**Theorem 2.8** ([26, Theorem 4.5]). *Consider two measures  $\mu, \nu$  on  $E$  with supports  $X$  and  $Y$ . Let  $\mu', \nu'$  be two measures with compact supports  $\Gamma$  and  $\Omega$  such that  $\Gamma \subset X$  and  $\Omega \subset Y$ . Then the interleaving distance  $d_i(V[X, d_\mu], V[Y, d_\nu])$  between the DTM-filtrations  $W[\mu]$  and  $W[\nu]$  is upper bounded by*

$$m^{-\frac{1}{2}} W_2(\mu, \mu') + m^{-\frac{1}{2}} W_2(\mu', \nu') + m^{-\frac{1}{2}} W_2(\nu', \nu) + c(\mu', m) + c(\nu', m),$$

where for any measure  $\tau$ , we define the quantity  $c(\tau, m) = \sup_{x \in \text{supp}(\tau)} d_{\tau, m}(x)$ .

Under a regularity assumption on  $\mu$ , we can restate Theorem 2.8 without mentioning the intermediate measures  $\mu'$  and  $\nu'$ . The proof is given in Appendix C.

**Corollary 2.9.** *Consider two probability measures  $\mu, \nu$  on  $E$ ,  $m \in (0, 1)$  and denote  $w = W_2(\mu, \nu)$ . Suppose that  $w \leq \frac{1}{4}$ , and that  $\mu$  satisfies the following for  $r < (\frac{m}{a})^{\frac{1}{d}}: \forall x \in \text{supp}(\mu), \mu(\mathcal{B}(x, r)) \geq ar^d$ . Then*

$$d_i(W[\mu], W[\nu]) \leq c_1 \left(\frac{w}{m}\right)^{\frac{1}{2}} + c'_1 m^{\frac{1}{d}},$$

with  $c_1 = 8\text{diam}(\text{supp}(\mu)) + 5$  and  $c'_1 = 2a^{-\frac{1}{d}}$ .

## 2.4 Model and hypotheses

**Model.** We consider an abstract  $\mathcal{C}^2$ -manifold  $\mathcal{M}_0$  of dimension  $d \geq 1$ ,  $E = \mathbb{R}^n$  the Euclidean space and a  $\mathcal{C}^2$ -immersion  $u: \mathcal{M}_0 \rightarrow E$ . We denote  $\mathcal{M} = u(\mathcal{M}_0)$ . Moreover, for any  $x_0 \in \mathcal{M}_0$ , we write  $x$  for  $u(x_0)$ ,  $T_{x_0}\mathcal{M}_0$  for the (abstract) tangent space of  $\mathcal{M}_0$  at  $x_0$ , and  $T_x\mathcal{M}$  for  $d_{x_0}u(T_{x_0}\mathcal{M}_0)$ , which is an affine subspace of  $E$ . Let  $\check{u}$  be the map

$$\check{u}: \mathcal{M}_0 \longrightarrow E \times \mathbf{M}(E) \\ x_0 \longmapsto \left( x, \frac{1}{d+2} p_{T_x\mathcal{M}} \right),$$

where  $p_{T_x\mathcal{M}}$  is the orthogonal projection matrix on  $T_x\mathcal{M}$ , and  $\mathbf{M}(E)$  the space of  $n \times n$  matrices. Note that the map  $\check{u}$  is a  $\mathcal{C}^1$ -immersion since  $u$  is  $\mathcal{C}^2$ . We define the *lifted manifold* as  $\check{\mathcal{M}} = \check{u}(\mathcal{M}_0)$ . We also consider a probability measure  $\mu_0$  on  $\mathcal{M}_0$ , and define  $\mu = u_*\mu_0$  and  $\check{\mu}_0 = \check{u}_*\mu_0$ . These several sets and measures fit in the following commutative diagrams:

$$\begin{array}{ccc} \mathcal{M}_0 & \xrightarrow{\check{u}} & \check{\mathcal{M}} \\ & \searrow u & \swarrow \text{proj} \\ & \mathcal{M} & \end{array} \qquad \begin{array}{ccc} \mu_0 & \xrightarrow{\check{u}_*} & \check{\mu}_0 \\ & \searrow u_* & \swarrow \text{proj}_* \\ & \mu & \end{array}$$

As explained in the introduction, the aim of this work is to estimate the homotopy type of  $\mathcal{M}_0$ , or its homology groups, from the measure  $\mu$ , or from a close measure  $\nu$ . We detail our method in Subsect. 5.1, and show that, by using DTM-filtrations, the problem boils down to estimating the measure  $\check{\mu}_0$  from  $\nu$ .

Besides, we endow  $\mathcal{M}_0$  with the Riemannian structure given by the immersion  $u$ . For every  $x_0 \in \mathcal{M}_0$ , the second fundamental form of  $\mathcal{M}_0$  at  $x_0$  is denoted  $\Pi_{x_0}$ , and the exponential map is denoted  $\exp_{x_0}^{\mathcal{M}_0}$ . We shall also consider the map  $\exp_x^{\mathcal{M}}: T_x\mathcal{M} \rightarrow \mathcal{M}$ , the *exponential map seen in  $\mathcal{M}$* , defined as  $u \circ \exp_{x_0}^{\mathcal{M}_0} \circ (d_{x_0}u)^{-1}$ .

**Notation conventions.** In the rest of this paper, symbols with 0 as a subscript shall refer to quantities associated to  $\mathcal{M}_0$ . For instance, a point of  $\mathcal{M}_0$  may be denoted  $x_0$ , and a curve on  $\mathcal{M}_0$  may be denoted  $\gamma_0$ . Symbols with a caron accent shall refer to quantities associated to  $\check{\mathcal{M}}$ , such as a point  $\check{x}$ , or a curve  $\check{\gamma}$ . Symbols with no such subscript or accent shall refer to quantities associated to  $\mathcal{M}$ , such as  $x$  or  $\gamma$ . In order to simplify the notations, we consider the following convention:

Dropping the 0 subscript to a symbol shall correspond to applying the map  $u$ .  
Dropping the 0 subscript to a symbol and adding a caron accent shall correspond to applying the map  $\check{u}$ .

For instance, if  $x_0$  is a point of  $\mathcal{M}_0$ , then  $x$  represents  $u(x_0)$ , and  $\check{x}$  represents  $\check{u}(x_0)$ . Similarly, if  $\gamma_0: I \rightarrow \mathcal{M}_0$  is a map, then  $\gamma$  represents  $u \circ \gamma_0$ , and  $\check{\gamma}$  represents  $\check{u} \circ \gamma_0$ . Note that it is possible to have  $x = y$  but  $T_x \mathcal{M} \neq T_y \mathcal{M}$ . When writing  $T_x \mathcal{M}$ , we will always refer to an implicit point  $x_0 \in \mathcal{M}_0$ .

**Hypotheses.** Throughout the paper, we shall refer to the four hypotheses listed below. The first one is a transversity-like condition.

**Hypothesis 1.** For every  $x_0, y_0 \in \mathcal{M}_0$  such that  $x_0 \neq y_0$  and  $x = y$ , we have  $T_x \mathcal{M} \neq T_y \mathcal{M}$ .

This Hypothesis 1 ensures that the  $\mathcal{C}^1$ -immersion  $\check{u}$  is injective, hence that it is a  $\mathcal{C}^1$ -diffeomorphism, since its domain  $\mathcal{M}_0$  is compact. As a consequence, the lifted manifold  $\check{\mathcal{M}}$  is a submanifold of  $E \times \mathbf{M}(E)$ , with the same homotopy type than  $\mathcal{M}_0$ . This allows to recover the homology of  $\mathcal{M}_0$  from  $\check{\mathcal{M}}$ .

**Hypothesis 2.** The operator norm of the second fundamental form of  $\mathcal{M}_0$  at each point is bounded by  $\rho > 0$ .

In Hypothesis 2, we consider that  $\mathcal{M}_0$  is endowed the Riemannian structure given by the immersion  $u$ . According to Equation (1), this hypothesis implies the following key property: if  $\gamma_0: I \rightarrow \mathcal{M}_0$  is an arc-length parametrized geodesic of class  $\mathcal{C}^2$ , then for all  $t \in I$ , we have  $\|\check{\gamma}(t)\| \leq \rho$ . In particular, we can use the Lemmas 2.3 and 2.4.

**Hypothesis 3.** The measure  $\mu_0$  admits a density  $f_0$  on  $\mathcal{M}_0$ . Moreover,  $f_0$  is  $L_0$ -Lipschitz (with respect to the geodesic distance) and bounded by  $f_{\min}, f_{\max} > 0$ .

In Hypothesis 3, we consider that  $\mathcal{M}_0$  is endowed with the volume measure  $\mathcal{H}_{\mathcal{M}_0}^d$ , that is, the measure obtained by pulling back the  $d$ -dimensional Hausdorff measure  $\mathcal{H}^d$  on  $E$  via the immersion  $u$ . Note that this may not be the uniform measure on  $\mathcal{M}$  (the volume is not renormalized). By assumption,  $\mu_0$  is a probability measure, hence the integral  $\int f_0 \cdot d\mathcal{H}_{\mathcal{M}_0}^d$  is equal to 1.

In order to state the fourth hypothesis, we need the notion of *normal reach*, that we will define in Subsect. 3.1. Roughly speaking, the normal reach is a map  $\lambda_0: \mathcal{M}_0 \rightarrow [0, +\infty)$  that indicates how close the point  $x_0$  is from a self-intersection. We remind the reader that we use the sublevel set notation  $\lambda_0^r = \lambda_0^{-1}([0, r])$ .

**Hypothesis 4.** There exists  $c_4 \geq 0$  and  $r_4 > 0$  such that, for every  $r \in [0, r_4)$ ,  $\mu_0(\lambda_0^r) \leq c_4 r$ .

This hypothesis will only be used in the last part of this paper, when gathering our results about tangent space estimation and stability of measures. Thanks to it, we will be able to subdivide  $\mathcal{M}_0$  in two sets that involve a different analysis: the points with small normal reach (points in  $\lambda_0^r$ ) and points with large normal reach (in  $\mathcal{M}_0 \setminus \lambda_0^r$ ).

The following table lists which hypotheses will be invoked in each subsection of the paper.

Subsection	3.1	3.2	3.3	4.2	4.3	4.4	4.5	5
Hyp. used	2	2, 3	1', 2, 3	2, 3	2, 3	2, 3	1, 2, 3, 4	1, 2, 3, 4

In Subsect. 4.3 and 4.4 we will introduce to a new set of hypotheses: 5, 6, and 7. We will show that these hypotheses are consequences of 2 and 3. However, referring to these new hypotheses will simplify the exposition, and will allow to state our results in a more general setting.

Concerning the naturality of the hypotheses, note that Hypothesis 1 is necessary to ensure that our problem is well-posed. To see this, consider the subset  $\mathcal{M} \subset \mathbb{R}^2$  consisting of two tangent circles. As depicted in Figure 7,  $\mathcal{M}$  may be the immersion of the abstract manifold  $\mathcal{M}_0$  being the disjoint union of two circles, or of  $\mathcal{M}'_0$  being a circle. In this case, the observation of  $\mathcal{M}$  cannot discriminate between  $\mathcal{M}_0$  and  $\mathcal{M}'_0$ . However, these immersions would not satisfy Hypothesis 1. Note that, in this example, the immersion  $\mathcal{M}'_0 \rightarrow \mathcal{M}$  is only  $\mathcal{C}^1$ , but one can easily design a similar example of regularity  $\mathcal{C}^2$ .

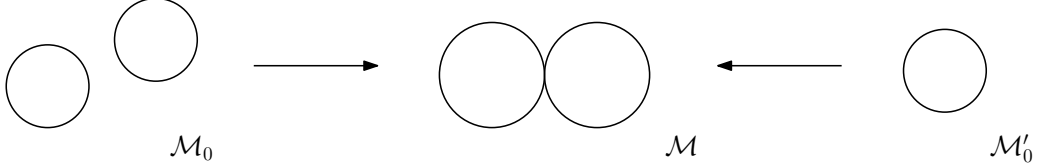


Figure 7: According to Hypothesis 1,  $\mathcal{M}$  cannot be two tangent circle.

In the literature, works often consider submanifolds, and Hypothesis 2 is usually stated as a lower bound on the reach  $\tau$ . Our hypothesis, stated as an upper bound on the norm  $\rho$  of the second fundamental forms of the immersion, is weaker. Indeed, under the assumption that the immersion is  $\mathcal{C}^2$ , we have  $\rho \leq \frac{1}{\tau}$ . The advantage of our formulation is that, in the case of an immersed manifold, the reach may be zero, hence cannot be used. Note that Hypothesis 2 is equivalent to the following property: there exists a function  $\alpha: [0, +\infty) \rightarrow [0, +\infty)$  such that  $\lim_{r \rightarrow 0} \alpha = 1$ , and such that for any  $x_0 \in \mathcal{M}_0$ , the image of the geodesic ball  $u(\mathcal{B}_{\mathcal{M}_0}(x_0, r))$  has reach lower bounded by  $\alpha(r) \frac{1}{\rho}$ . The fact that Hypothesis 2 implies this statement is a consequence of the proof of [2, Proposition 6.1], and the converse is a consequence of the proof of Proposition 3.8.

The introduction of constants in Hypothesis 3 will allow derive explicit bounds for our method. Note that we do not suppose that the measure  $\mu$  is given as an input, but only a close measure  $\nu$  with respect to the Wasserstein distance. There is no hypothesis concerning the measure  $\nu$ .

Last, we think that Hypothesis 4 is a consequence of Hypotheses 1, 2 and 3, but we have not been able to prove it yet. As a partial result, we prove that it is a consequence of Hypotheses 1', 2 and 3, where Hypothesis 1' is a strenghtening of Hypothesis 1 (see Proposition 3.19). We also show that Hypothesis 4 is a consequence of Hypotheses 1, 2, 3 and  $\dim(\mathcal{M}_0) = 1$  (see Remark 3.21).

### 3 Reach of an immersed manifold

In this section, we introduce a new notion of reach, adapted to the immersed manifolds, and derive technical results that will be useful in the rest of the paper. As an introduction, we consider an embedded manifold  $u: \mathcal{M}_0 \rightarrow \mathcal{M} \subset \mathbb{R}^n$  with positive reach  $\tau$ . Let  $x_0, y_0$  be two points of  $\mathcal{M}_0$ . We wish to compare their geodesic distance  $d_{\mathcal{M}_0}(x_0, y_0)$  and their Euclidean distance  $\|y - x\|$ . A first inequality is true in general:

$$\|y - x\| \leq d_{\mathcal{M}_0}(x_0, y_0).$$

Moreover, if they are close enough in geodesic distance—say  $d_{\mathcal{M}_0}(x_0, y_0) \leq \tau$  for instance—then the inequality  $\rho \leq \frac{1}{\tau}$  and Lemma 2.3 Point 3 yields

$$d_{\mathcal{M}_0}(x_0, y_0) \leq 2\|x - y\|.$$

This section is devoted to obtaining such a converse inequality when the manifold  $\mathcal{M}_0$  is only immersed, not embedded. In this case, the condition  $d_{\mathcal{M}_0}(x_0, y_0) \leq \frac{1}{\tau}$  has to be turned into an upper bound on  $d_{\mathcal{M}_0}(x_0, y_0)$  that depends on  $x_0$  and  $y_0$  (as we will obtain in Lemma 3.4).

### 3.1 Normal reach

We consider an immersion  $u: \mathcal{M}_0 \rightarrow \mathcal{M} \subset E$  which satisfies Hypothesis 2.

**Definition 3.1.** For every  $x_0 \in \mathcal{M}_0$ , let  $\Lambda(x_0) = \{y_0 \in \mathcal{M}_0 \mid y_0 \neq x_0, x - y \perp T_{y_0}\mathcal{M}\}$ . The *normal reach* of  $\mathcal{M}_0$  at  $x_0$  is defined as:

$$\lambda_0(x_0) = \inf_{y_0 \in \Lambda(x_0)} \|x - y\|.$$

Observe that if  $x_0, y_0$  are distinct points of  $\mathcal{M}_0$  with  $x = y$ , then  $x - y$  is orthogonal to any vector, hence  $\lambda_0(x_0) = \|x - y\| = 0$ . Hence we can define the *normal reach seen in  $\mathcal{M}$* , denoted  $\lambda: \mathcal{M} \rightarrow \mathbb{R}$ , as

$$\lambda(x) = \begin{cases} \lambda_0(u^{-1}(x)) & \text{if } x \text{ has only one preimage,} \\ 0 & \text{else.} \end{cases}$$

It satisfies the relation  $\lambda_0 = \lambda \circ u$ .

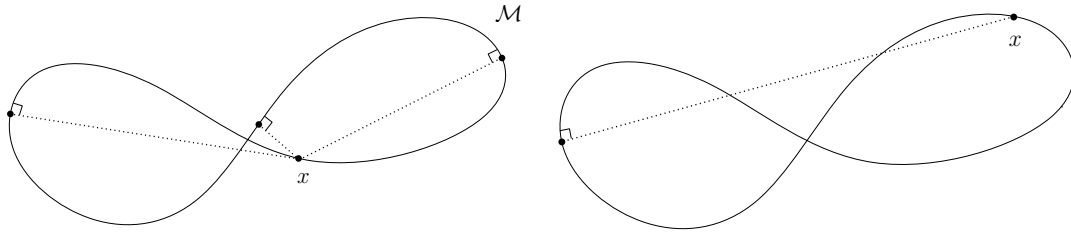


Figure 8: The set  $\Lambda(x_0)$  from Definition 3.1, for two different points  $x_0$ .

Note that  $\Lambda(x_0)$  is closed, hence the infimum of Definition 3.1 is attained. Indeed, we can write  $\Lambda(x_0) = L \setminus \{x_0\}$ , with  $L = \{y_0 \in \mathcal{M}_0 \mid x - y \perp T_{y_0}\mathcal{M}\}$ . The set  $L$  is closed since it is the preimage of  $\{0\}$  by the continuous map  $y_0 \mapsto \|p_{T_{y_0}\mathcal{M}}(x - y)\|$ . Furthermore,  $\{x_0\}$  is an isolated point of  $\Lambda(x_0)$ , since Lemma 2.3 Point 5 says that, for every  $y_0$  in the geodesic ball  $\mathcal{B}_{\mathcal{M}_0}(x_0, \frac{1}{\rho})$  such that  $y_0 \neq x_0$ , the vector  $x - y$  is not orthogonal to  $T_{y_0}\mathcal{M}$ , hence  $y_0 \notin L$ .

**Example 3.1.** Suppose that  $\mathcal{M}$  is the lemniscate of Bernoulli, with diameter 2. Figure 9 represents the values of the normal reach  $\lambda: \mathcal{M} \rightarrow \mathbb{R}$ . Observe that  $\lambda$  is not continuous.

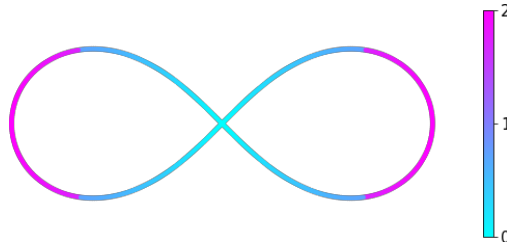


Figure 9: Values of the normal reach on the lemniscate of Bernoulli.

**Remark 3.2.** The normal reach  $\lambda_0$  is *lower semi-continuous*, that is, for any sequence  $(x_0^n)_{n \geq 0}$  of  $\mathcal{M}_0$  converging to a point  $x_0 \in \mathcal{M}_0$ , we have  $\liminf_{n \rightarrow \infty} \lambda_0(x_0^n) \geq \lambda_0(x_0)$ . Indeed, by definition, we can choose for every  $n \geq 0$  a point  $y_0^n \in \mathcal{M}_0$  such that  $x_0^n \neq y_0^n$ ,  $x^n - y^n \perp T_{y_0^n}\mathcal{M}$  and  $\lambda_0(x_0^n) = \|x^n - y^n\|$ . Moreover, since  $\mathcal{M}_0$  is compact, the sequence  $(y_0^n)_{n \geq 0}$  admits an accumulation point  $y_0$  such that  $\liminf_{n \rightarrow \infty} \lambda_0(x_0^n) = \|x - y\|$ . By continuity, we have  $x - y \perp T_{y_0}\mathcal{M}$ . Moreover, since  $d_{\mathcal{M}_0}(x_0^n, y_0^n) \geq \frac{1}{\rho}$  for all  $n \geq 0$  by Lemma 2.3 Point 5, we deduce that  $x_0 \neq y_0$ . Consequently, we have  $y_0 \in \Lambda(x_0)$ , and  $\lambda_0(x_0) \leq \|x - y\|$ .

Here is a key property of the normal reach:

**Lemma 3.3.** Let  $x_0 \in \mathcal{M}_0$ . Let  $r \geq 0$  such that  $r < \lambda(x)$ . Then  $u^{-1}(\overline{\mathcal{B}}(x, r))$  is connected.



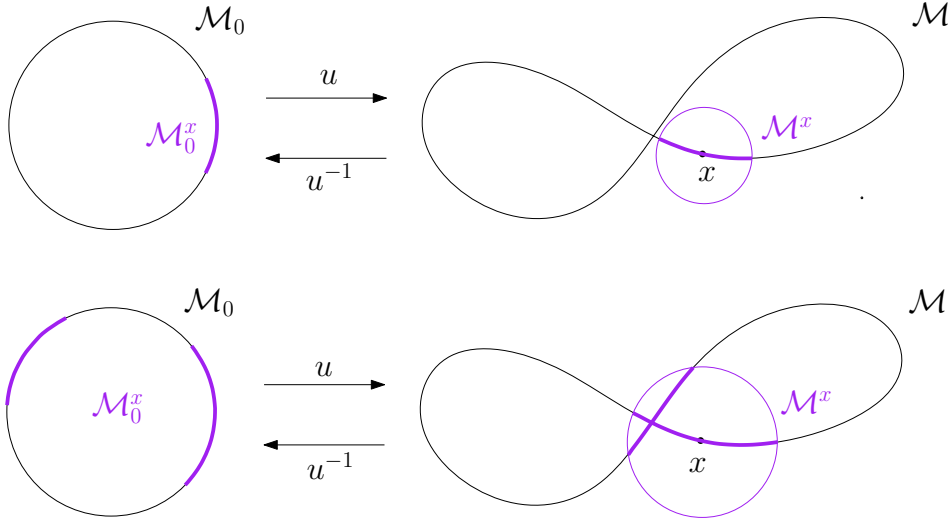


Figure 10: Top: the set  $u^{-1}(\overline{\mathcal{B}}(x, r))$ , with  $r < \lambda(x)$ , is connected. Bottom: when  $r \geq \lambda(x)$ , it may not be connected.

*Proof.* Write  $\mathcal{M}^x = \mathcal{M} \cap \overline{\mathcal{B}}(x, r)$  and  $\mathcal{M}_0^x = u^{-1}(\mathcal{M}^x)$ . By contradiction, suppose that  $\mathcal{M}_0^x$  is not connected. Let  $C \subset \mathcal{M}_0^x$  be a connected component which does not contain  $x_0$ . Since  $C$  is compact, we can consider a minimizer  $y_0$  of  $\{\|x - y\| \mid y_0 \in C\}$ . It is clear that  $y$  satisfies  $x - y \perp T_y \mathcal{M}$ , otherwise it would not be a local minimizer. Now, the properties  $x - y \perp T_y \mathcal{M}$  and  $x_0 \neq y_0$  imply that  $\|x - y\| \geq \lambda(x)$ , which contradicts  $r < \lambda(x)$ .  $\square$

The following lemma is the counterpart of [2, Proposition 6.3] for the normal reach. It allows to compare the geodesic and Euclidean distances by only imposing a condition on the last one.

**Lemma 3.4.** *Let  $x_0, y_0 \in \mathcal{M}_0$ . Denote  $r = \|x - y\|$  and  $\delta = d_{\mathcal{M}_0}(x_0, y_0)$ . Suppose that  $r < \min\left(\frac{1}{2\rho}, \lambda(x)\right)$ . Then*

$$\delta \leq c_5(\rho r)r \quad \text{where} \quad c_5(t) = \frac{1}{t} \left(1 - \sqrt{1 - 2t}\right).$$

*In other words, the following inclusion holds:  $u^{-1}(\overline{\mathcal{B}}(x, r)) \subset \overline{\mathcal{B}}_{\mathcal{M}_0}(x_0, c_5(\rho r)r)$ .*

Note that, for  $t < \frac{1}{2}$ , we have the inequalities  $1 \leq c_5(t) \leq 1 + 2t < 2$ .

*Proof.* Denote  $\mathcal{M}^x = \mathcal{M} \cap \overline{\mathcal{B}}(x, r)$  and  $\mathcal{M}_0^x = u^{-1}(\mathcal{M}^x)$ . Let  $C_0^x$  be the connected component of  $x_0$  in  $\mathcal{M}_0^x$ . Let us show that  $C_0^x$  is included in the closed geodesic ball  $\overline{\mathcal{B}}_{\mathcal{M}_0}(x_0, c_5(\rho r)r)$ .

Let  $\varepsilon > 0$  be such that  $r < r + \varepsilon < \frac{1}{2\rho}$ . For any tangent vector  $v_0 \in T_{x_0} \mathcal{M}_0$  of unit norm, consider the unit-speed geodesic  $\gamma_0: I \rightarrow \mathcal{M}_0$  with  $\gamma(0) = x_0$  and  $\dot{\gamma}(0) = v_0$ . According to Lemma 2.4 Points 2 and 3, for any  $t$  in  $\left(\frac{1}{\rho}(1 - \sqrt{1 - 2\rho r}), \frac{1}{\rho}(1 + \sqrt{1 - 2\rho r})\right)$ , we have  $\|x - \gamma(t)\| > r$ . Consequently,  $C_0^x$  and the geodesic sphere  $\partial \mathcal{B}_{\mathcal{M}_0}(x_0, t)$  are disjoint. We deduce that  $C_0^x \subset \mathcal{B}_{\mathcal{M}_0}(x_0, t)$ . In particular, we have  $C_0^x \subset \overline{\mathcal{B}}_{\mathcal{M}_0}(x_0, t^*)$ , where

$$t^* = \frac{1}{\rho}(1 - \sqrt{1 - 2\rho r}) = c_5(\rho r)r.$$

Besides,  $\mathcal{M}_0^x$  is connected by Lemma 3.3. We deduce that  $\mathcal{M}_0^x \subset \overline{\mathcal{B}}_{\mathcal{M}_0}(x_0, c_5(\rho r)r)$ . In particular,  $y_0$  must satisfy  $d_{\mathcal{M}_0}(x_0, y_0) \leq c_5(\rho r)r$ , and we deduce the result.  $\square$

Following the same idea, we can prove the following lemma. It states that normal reach  $\lambda_0(x_0)$  can be understood as the minimal distance  $\|x - y\|$  for points  $y_0$  far enough from  $x_0$  in  $\mathcal{M}_0$

**Lemma 3.5.** *Let  $x_0, y_0 \in \mathcal{M}_0$  such that  $\|x - y\| < \frac{1}{2\rho}$  and  $d_{\mathcal{M}_0}(x_0, y_0) \geq 4\|x - y\|$ . Then there exists a  $z_0 \in \mathcal{M}_0$  such that  $d_{\mathcal{M}_0}(y_0, z_0) < 2\|x - y\|$ ,  $\|x - z\| \leq \|x - y\|$  and  $x - z \perp T_{z_0}\mathcal{M}$ . Consequently,  $\lambda_0(x_0) \leq \|x - y\|$ .*

*Proof.* Denote  $r = \|x - y\|$  and  $\mathcal{M}_0^x = u^{-1}(\overline{\mathcal{B}}(x, r))$ . Let  $C_0^x$  denote the connected component of  $x_0$  in  $\mathcal{M}_0^x$ . As we have seen in the proof of Lemma 3.4,  $C_0^x$  is included in the closed geodesic ball  $\overline{\mathcal{B}}_{\mathcal{M}_0}(x_0, c_5(\rho r)r)$ . Since  $r < \frac{1}{2\rho}$ , we have the inequality  $c_5(\rho r) < 2$ , and we deduce  $C_0^x \subset \mathcal{B}_{\mathcal{M}_0}(x_0, 2r)$ . Now, let  $C_0^y$  denote the connected component of  $y_0$  in  $\mathcal{M}_0^x$ . Similarly, we have  $C_0^y \subset \mathcal{B}_{\mathcal{M}_0}(x_0, 2r)$ . Since  $d_{\mathcal{M}_0}(x_0, y_0) \geq 4r$ , we deduce that  $C_0^x$  and  $C_0^y$  are disjoint. Now, let  $z_0$  be a minimizer of  $z_0 \mapsto \|x - z\|$  on  $C_0^y$ . We have  $z_0 \neq x_0$ . Moreover,  $x - z \perp T_{z_0}\mathcal{M}$ , otherwise  $z$  would not be a minimizer. Hence, by definition of the normal reach,  $\lambda_0(x_0) \leq \|x - z\| \leq \|x - y\|$ .  $\square$

The following proposition connects the normal reach to the usual notion of reach, in the case where  $\mathcal{M}_0$  is embedded.

**Proposition 3.6.** *Suppose that  $u: \mathcal{M}_0 \rightarrow \mathcal{M} \subset E$  is a  $\mathcal{C}^2$ -embedding. Let  $\tau > 0$  be the reach of  $\mathcal{M}$ . We have*

$$\tau = \min\left(\frac{1}{\rho_*}, \frac{1}{2}\lambda_*\right),$$

where  $\rho_*$  is the supremum of the operator norm of the second fundamental forms of  $\mathcal{M}_0$ , and  $\lambda_* = \inf_{x \in \mathcal{M}} \lambda(x)$  is the infimum of the normal reach.

*Proof.* We first prove that  $\tau \geq \min\left(\frac{1}{\rho_*}, \frac{1}{2}\lambda_*\right)$ . According to Theorem 2.2, two cases may occur: the reach is either caused by a bottleneck or by curvature. In the first case, there exists  $x, y \in \mathcal{M}$  and  $z \in \text{med}(\mathcal{M})$  with  $\|x - y\| = 2\tau$  and  $\|x - z\| = \|y - z\| = \tau$ . We deduce that  $x - y \perp T_y\mathcal{M}$ . Hence by definition of  $\lambda(x)$ ,

$$\lambda(x) \leq \|x - y\| = 2\|x - z\| \leq 2\tau.$$

In the second case, there exists  $x \in \mathcal{M}$  and an arc-length parametrized geodesic  $\gamma: I \rightarrow \mathcal{M}$  such that  $\gamma(0) = x$  and  $\|\dot{\gamma}(0)\| = \frac{1}{\tau}$ . But  $\|\dot{\gamma}(0)\| \leq \rho_*$ , hence  $\frac{1}{\tau} \leq \rho_*$ . This disjunction shows that  $\tau \geq \min\left(\frac{1}{\rho_*}, \frac{1}{2}\lambda_{\min}\right)$ .

We now prove that  $\tau \leq \min\left(\frac{1}{\rho_*}, \frac{1}{2}\lambda_*\right)$ . The inequality  $\tau \leq \frac{1}{\rho_*}$  appears in [2, Proposition 6.1]. To prove  $\tau \leq \frac{1}{2}\lambda_*$ , consider any  $x_0 \in \mathcal{M}_0$ . Let  $y_0 \in \Lambda(x_0)$  such that  $\|x - y\|$  is minimal. Using Theorem 2.1 and the property  $x - y \perp T_{y_0}\mathcal{M}$ , we immediately have

$$\tau \leq \frac{\|x - y\|^2}{2\text{dist}(y - x, T_{y_0}\mathcal{M})} = \frac{\|x - y\|}{2} = \frac{\lambda(x)}{2},$$

and the result follows.  $\square$

**Remark 3.7.** We can generalize Proposition 3.6 as follows: suppose that  $u: \mathcal{M}_0 \rightarrow \mathcal{M} \subset E$  is a  $\mathcal{C}^2$ -immersion (potentially an embedding) that satisfies Hypothesis 1. Let  $\tau \geq 0$  be the reach of  $\mathcal{M}$ . We have

$$\tau = \min\left(\frac{1}{\rho_*}, \frac{1}{2}\lambda_*\right).$$

In fact, if  $u$  is not an embedding, we can show that  $\tau = 0$  and  $\lambda_* = 0$ . Indeed, if  $x \in \mathcal{M}$  is a point that admits several preimages by  $u$ , we have seen that  $\lambda(x) = 0$ . On the other hand, by Hypothesis 1, the tangent cone  $\text{Tan}(\mathcal{M}, x)$  is not an affine subspace but an union of several affine subspaces. In this case, we see that Theorem 2.1 cannot hold, hence that  $\tau = 0$ .

Note that if Hypothesis 1 is not satisfied, it is possible to have  $\tau > 0$  but  $\lambda_* = 0$ . This would be the case for any non-injective immersion  $u: \mathcal{M}_0 \rightarrow \mathcal{M}$  such that its image  $\mathcal{M}$  is a  $\mathcal{C}^2$ -submanifold.

As shown by the previous remark, when  $\mathcal{M}$  is not a submanifold, global quantities such as the reach  $\tau$  or the minimal normal reach  $\lambda_*$  are zero. However, as shown by the following proposition, the normal reach gives a scale at which  $\mathcal{M}$  still behaves well. Note that we shall not make use of this result in the rest of the paper.

**Proposition 3.8.** *Assume that  $\mathcal{M}_0$  satisfies Hypothesis 2. Let  $x \in \mathcal{M}_0$  and  $r < \min\left(\frac{1}{4\rho}, \lambda(x)\right)$ . Then  $\mathcal{M} \cap \overline{\mathcal{B}}(x, r)$  is a set of reach at least  $\frac{1-2\rho r}{\rho}$ .*

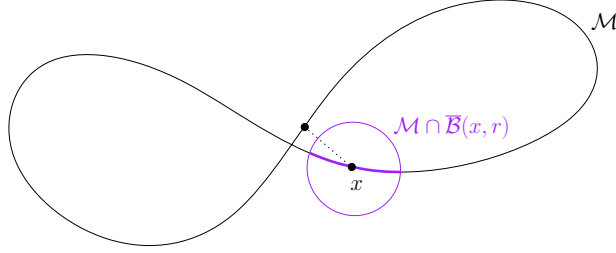


Figure 11: The set  $\mathcal{M} \cap \overline{\mathcal{B}}(x, r)$  has positive reach.

*Proof.* Denote  $\mathcal{M}^x = \mathcal{M} \cap \overline{\mathcal{B}}(x, r)$  and  $\mathcal{M}_0^x = u^{-1}(\mathcal{M}^x)$ . In order to give a bound on the reach of  $\mathcal{M}^x$ , we will use the characterization of Theorem 2.1. First, let us prove that for every  $y_0, z_0 \in \mathcal{M}_0^x$ ,

$$\text{dist}(z - y, T_y \mathcal{M}) \leq \frac{\rho}{2(1-2\rho r)} \|z - y\|^2.$$

Let  $y_0, z_0 \in \mathcal{M}_0^x$ , and  $\delta = d_{\mathcal{M}_0}(y_0, z_0)$ . Lemma 2.3 Point 3 gives  $\delta \leq \frac{1}{1-\frac{\rho}{2}\delta} \|y - z\|$ . Moreover,  $\delta \leq d_{\mathcal{M}_0}(y_0, x_0) + d_{\mathcal{M}_0}(x_0, z_0) \leq 2c_5(\rho r)r$ . Hence,

$$\frac{1}{1-\frac{\rho}{2}\delta} \leq \frac{1}{1-c_5(\rho r)\rho r} = \frac{1}{\sqrt{1-2\rho r}},$$

and we deduce that

$$\delta \leq \frac{1}{\sqrt{1-2\rho r}} \|y - z\|. \quad (5)$$

Besides, Lemma 2.3 Point 2 gives  $\text{dist}(z - y, T_y \mathcal{M}) \leq \frac{\rho}{2} \delta^2$ , and combining these two inequalities yields  $\text{dist}(z - y, T_y \mathcal{M}) \leq \frac{\rho}{2(1-2\rho r)} \|z - y\|^2$ .

Secondly, let us prove that

$$\text{dist}(z - y, \text{Tan}(\mathcal{M}^x, y)) \leq \frac{\rho}{2(1-2\rho r)} \|z - y\|^2, \quad (6)$$

where  $\text{Tan}(\mathcal{M}^x, y)$  is the tangent cone at  $y$  of the closed set  $\mathcal{M}^x$ . According to Equation (5), it is enough to prove that  $\text{Tan}(\mathcal{M}^x, y) = T_y \mathcal{M}$ . We shall prove that  $\mathcal{M}$  does not self-intersect in  $\overline{\mathcal{B}}(x, r)$ . According to Lemma 3.4, we have  $\mathcal{M}_0^x \subset \mathcal{B}_{\mathcal{M}_0}\left(x_0, \frac{1}{\rho}\right)$ . Using Lemma 2.3 Point 4, we get that  $u$  is injective on  $\mathcal{M}_0^x$ , as wanted. To conclude the proof, it follows from Theorem 2.1 and Equation (6) that  $\mathcal{M}^x$  has reach at least  $\frac{1-2\rho r}{\rho}$ .  $\square$

### 3.2 Probabilistic bounds under normal reach conditions

We now consider  $\mathcal{M}_0$  and  $\mu_0$  which satisfy Hypotheses 2 and 3. The aim of this subsection is to provide a quantitative control of the measure  $\mu = u_* \mu_0$ , that is, bounds on the measure of balls and annuli (see Propositions 3.13 and 3.14). We do so by pulling-back  $\mu$  on the tangent spaces  $T_x \mathcal{M}$ , where it is simpler to compute integrals (see Lemma 3.11).

Recall that the exponential map of  $\mathcal{M}_0$  at a point  $x_0 \in \mathcal{M}_0$  is denoted

$$\exp_{x_0}^{\mathcal{M}_0} : T_{x_0}\mathcal{M}_0 \rightarrow \mathcal{M}_0.$$

To ease the reading of this subsection, we introduce the *exponential map seen in  $\mathcal{M}$* , denoted  $\exp_x^{\mathcal{M}} : T_x\mathcal{M} \rightarrow \mathcal{M}$ . It is defined as

$$\exp_x^{\mathcal{M}} = u \circ \exp_{x_0}^{\mathcal{M}_0} \circ (d_{x_0}u)^{-1}.$$

Note that the map  $\exp_x^{\mathcal{M}}$  is well-defined, even if  $x$  is a self-intersection point of  $\mathcal{M}$ . Indeed,  $\exp_x^{\mathcal{M}}$  will always refer implicitly to a choice of point  $x_0$  such that  $x = u(x_0)$ . This is consistent with the notation conventions of Subsect. 2.4. This map fits in the following commutative diagram:

$$\begin{array}{ccc} \mathcal{M}_0 & \xrightarrow{u} & \mathcal{M} \\ \exp_{x_0}^{\mathcal{M}_0} \uparrow & & \uparrow \exp_x^{\mathcal{M}} \\ T_{x_0}\mathcal{M}_0 & \xrightarrow{d_{x_0}u} & T_x\mathcal{M} \end{array}$$

We also define the map  $\overline{\exp}_x^{\mathcal{M}}$  as the restriction of  $\exp_x^{\mathcal{M}}$  to the closed ball  $\overline{\mathcal{B}}_{T_x\mathcal{M}}\left(0, \frac{2}{\rho}\right)$  of  $T_x\mathcal{M}$ . It is injective by Lemma 2.3 Point 4, and its image is  $u\left(\overline{\mathcal{B}}_{\mathcal{M}_0}\left(x_0, \frac{2}{\rho}\right)\right)$ . Moreover, for any  $r \leq \min\left(\frac{1}{2\rho}, \lambda(x)\right)$ ,  $u\left(\overline{\mathcal{B}}_{\mathcal{M}_0}\left(x_0, \frac{2}{\rho}\right)\right)$  contains  $\mathcal{M} \cap \overline{\mathcal{B}}(x, r)$  by Lemma 3.4, hence we can consider its inverse

$$(\overline{\exp}_x^{\mathcal{M}})^{-1} : \mathcal{M} \cap \overline{\mathcal{B}}(x, r) \longrightarrow T_x\mathcal{M}. \quad (7)$$

The next lemma gathers previous results. We remind the reader that the  $d$ -dimensional Jacobian has been defined in Subsect. 2.1.

**Lemma 3.9.** *Let  $x_0 \in \mathcal{M}_0$  and  $r < \min\left(\frac{1}{2\rho}, \lambda(x)\right)$ . Denote  $\overline{\mathcal{B}}_0 = (\overline{\exp}_x^{\mathcal{M}})^{-1}(\overline{\mathcal{B}}(x, r))$ . We have the inclusions*

$$\overline{\mathcal{B}}_{T_x\mathcal{M}}(0, r) \subset \overline{\mathcal{B}}_0 \subset \overline{\mathcal{B}}_{T_x\mathcal{M}}(0, c_5(\rho r)).$$

Moreover, for all  $v \in \overline{\mathcal{B}}_0$ , the Jacobian  $J_v$  of  $\overline{\exp}_x^{\mathcal{M}}$ , is bounded by

$$\left(1 - \frac{(r\rho)^2}{6}\right)^d \leq J_v \leq (1 + (r\rho)^2)^d,$$

and these terms are bounded by  $J_{\min} = \left(\frac{23}{24}\right)^d$  and  $J_{\max} = \left(\frac{5}{4}\right)^d$ .

*Proof.* The inclusions come from Lemma 3.4. The bounds on the Jacobian come from Lemma 2.5 and the fact that  $c_5(\rho r)r \leq 2r \leq \frac{1}{\rho} \leq \frac{\pi}{2\sqrt{2}\rho}$  when  $r < \frac{1}{2\rho}$ .  $\square$

We now study the measure  $\mu$ . By definition, it is the push-forward of  $\mu_0$  by  $u$ . By applying the coarea formula, and in particular Equation (2), we obtain that  $\mu$  admits the following density against  $\mathcal{H}_{\mathcal{M}}^d$ , the  $d$ -dimensional Hausdorff measure restricted to  $\mathcal{M}$ :

$$f(x) = \sum_{x_0 \in u^{-1}(\{x\})} f_0(x_0).$$

Indeed, in this case, the Jacobian is always 1, since  $\mathcal{M}_0$  has been given the pull-back Riemannian metric. Note that if  $x$  has only one preimage by  $u$ —i.e., if  $\lambda(x) > 0$ —then  $f(x) = f_0 \circ u^{-1}(x)$ . In the rest of the paper, we will only use  $f$  on points  $x$  such that  $\lambda(x) > 0$ . This is motivated by the fact that in Sect. 4 and 5, we will assume Hypothesis 4, which gives that the measure of the set  $\{x_0 \in \mathcal{M}_0 \mid \lambda_0(x_0) = 0\}$  is zero. Moreover, we have a Lipschitz-like property for the density  $f$ , valid as long as the points are chosen far enough from the self-intersection of  $\mathcal{M}$ :

**Lemma 3.10.** For all  $x_0, y_0 \in \mathcal{M}_0$  such that  $\|x - y\| < \min\left(\frac{1}{2\rho}, \lambda(x)\right)$ , we have

$$|f(x) - f(y)| \leq 2L_0 \|x - y\|.$$

*Proof.* Recall that, by Hypothesis 3, the density  $f_0$  is  $L_0$ -Lipschitz with respect to the geodesic distance: for all  $x_0, y_0 \in \mathcal{M}_0$ , we have  $|f_0(x_0) - f_0(y_0)| \leq L_0 \cdot d_{\mathcal{M}_0}(x_0, y_0)$ . To prove the lemma, we start with the case where  $y$  has only one preimage by  $u$ , so that we can write  $f(y) = f_0 \circ u^{-1}(y)$ . Since  $\|x - y\| < \lambda(x)$  by assumption, we have  $0 < \lambda(x)$ , hence  $x$  also has only one preimage. Now we have

$$\begin{aligned} |f(x) - f(y)| &= |f_0 \circ u^{-1}(x) - f_0 \circ u^{-1}(y)| \\ &\leq L_0 \cdot d_{\mathcal{M}_0}(u^{-1}(x), u^{-1}(y)) \\ &\leq 2L_0 \|x - y\|, \end{aligned}$$

where we used Lemma 3.4 on the last inequality. Now we prove that  $\|x - y\| < \min\left(\frac{1}{2\rho}, \lambda(x)\right)$  implies that  $y$  has only one preimage. Let  $r = \|x - y\|$ , and suppose by contradiction that  $y_0, y'_0$  are two distinct preimages. According to Lemma 2.3 Point 4,  $d_{\mathcal{M}_0}(y_0, y'_0) \geq \frac{2}{\rho}$ . But Lemma 3.4 says that  $u^{-1}(\mathcal{B}(x, r)) \subset \mathcal{B}_{\mathcal{M}_0}(x_0, 2r) \subset \mathcal{B}_{\mathcal{M}_0}\left(x_0, \frac{1}{\rho}\right)$ , which yields the contradiction  $d_{\mathcal{M}_0}(y_0, y'_0) \leq d_{\mathcal{M}_0}(y_0, x_0) + d_{\mathcal{M}_0}(y'_0, x_0) < \frac{1}{\rho}$ .  $\square$

We now state the key lemma of this subsection, that allows to go from a measure on  $\mathcal{M}_0$  to a measure on  $\mathcal{M}$ .

**Lemma 3.11.** Let  $x_0 \in \mathcal{M}_0$  and  $r < \min\left(\frac{1}{2\rho}, \lambda(x)\right)$ . Consider  $\mu_x$ , the measure  $\mu$  restricted to  $\overline{\mathcal{B}}(x, r)$ , define  $\overline{\mathcal{B}}_0 = (\overline{\text{exp}}_x^{\mathcal{M}})^{-1}(\overline{\mathcal{B}}(x, r))$  and the push-forward

$$\nu_x = \left(\overline{\text{exp}}_x^{\mathcal{M}}\right)_*^{-1} \mu_x,$$

where  $(\overline{\text{exp}}_x^{\mathcal{M}})^{-1}$  has been defined in Equation (7). The measure  $\nu_x$  admits the following density against the  $d$ -dimensional Hausdorff measure on  $T_x\mathcal{M}$ :

$$g(v) = f(\overline{\text{exp}}_x^{\mathcal{M}}(v)) \cdot J_v \cdot 1_{\overline{\mathcal{B}}_0}(v).$$

Moreover, for all  $v \in \overline{\mathcal{B}}_0$ , the map  $g$  satisfies  $|g(v) - g(0)| \leq c_7 r$ , where  $c_7 = 4L_0 J_{\max} + \frac{d}{2}\rho f_{\max}$ .

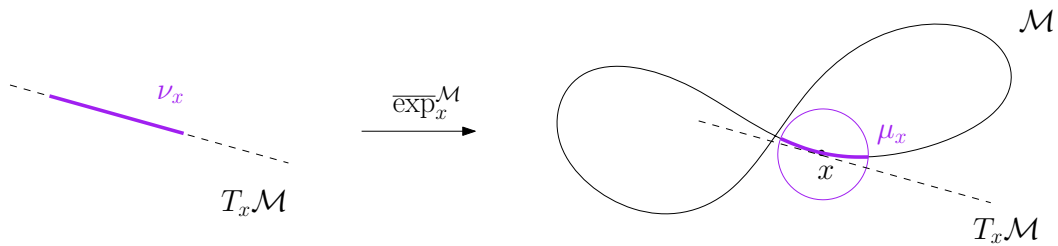


Figure 12: Measures involved in Lemma 3.11.

*Proof.* The expression of  $g$  comes from the coarea formula (Equation (2)) applied to the map  $(\overline{\text{exp}}_x^{\mathcal{M}})^{-1}: \mathcal{M} \cap \overline{\mathcal{B}}(x, r) \rightarrow \overline{\mathcal{B}}_0$ , and the measure  $\nu_x = (\overline{\text{exp}}_x^{\mathcal{M}})_*^{-1} \mu_x$ . To prove the inequality, observe that we can decompose

$$\begin{aligned} g(v) - g(0) &= f(\overline{\text{exp}}_x^{\mathcal{M}}(v)) J_v - f(\overline{\text{exp}}_x^{\mathcal{M}}(0)) J_0 \\ &= \left[ f(\overline{\text{exp}}_x^{\mathcal{M}}(v)) - f(\overline{\text{exp}}_x^{\mathcal{M}}(0)) \right] J_v + (J_v - J_0) f(\overline{\text{exp}}_x^{\mathcal{M}}(0)). \end{aligned}$$

On the one hand, using Lemma 3.10, we get

$$\begin{aligned} \left| f\left(\overline{\exp}_x^{\mathcal{M}}(v)\right) - f\left(\overline{\exp}_x^{\mathcal{M}}(0)\right) \right| &\leq 2L_0 \left\| \overline{\exp}_x^{\mathcal{M}}(v) - \overline{\exp}_x^{\mathcal{M}}(0) \right\| \\ &= 2L_0 \left\| u \circ \exp_{x_0}^{\mathcal{M}_0}(v) - u \circ \exp_{x_0}^{\mathcal{M}_0}(0) \right\| \\ &\leq 2L_0 \cdot d_{\mathcal{M}_0}(\overline{\exp}_{x_0}^{\mathcal{M}_0}(v), x_0) = 2L_0 \|v\|. \end{aligned}$$

On the other hand,  $J_0 = 1$  and  $\left(1 - \frac{(rp)^2}{6}\right)^d \leq J_v \leq (1 + (rp)^2)^d$  yield  $|J_v - J_0| \leq d(rp)^2 \leq \frac{d}{2}\rho r$ . Using the triangle inequality we see that

$$|g(v) - g(0)| \leq 2L_0 \|v\| J_{\max} + f_{\max} \frac{d}{2} \rho r \leq \left(4L_0 J_{\max} + f_{\max} \frac{d}{2} \rho\right) r,$$

as wanted.  $\square$

*Remark 3.12.* In the same vein as Lemma 3.11, define  $\overline{\exp}_{x_0}^{\mathcal{M}_0}$  to be the map  $\exp_{x_0}^{\mathcal{M}_0}$  restricted to  $\overline{\mathcal{B}}_{T_{x_0}\mathcal{M}_0}\left(0, \frac{2}{\rho}\right)$ . For any  $x_0 \in \mathcal{M}_0$ , let  $\mu_0^{x_0}$  be the measure  $\mu_0$  restricted to  $\overline{\mathcal{B}}_{\mathcal{M}_0}\left(x_0, \frac{2}{\rho}\right)$ , and define the measure

$$v_0 = (\overline{\exp}_{x_0}^{\mathcal{M}_0})^{-1} \mu_0^{x_0}.$$

Using the area formula, one shows that  $v_0$  admits the following density over the  $d$ -dimensional Hausdorff measure on  $T_{x_0}\mathcal{M}_0$ :

$$g_0(v) = f_0\left(\overline{\exp}_{x_0}^{\mathcal{M}_0}(v)\right) \cdot J_v \cdot 1_{\overline{\mathcal{B}}_{T_{x_0}\mathcal{M}_0}\left(0, \frac{2}{\rho}\right)}(v).$$

Now we can use the density  $g$  of Lemma 3.11 to derive explicit bounds on  $\mu$ . We remind the reader that  $V_d$  denote the volume of the unit ball of  $\mathbb{R}^d$ .

**Proposition 3.13.** *Let  $x_0 \in \mathcal{M}_0$ ,  $r \leq \min\left(\frac{1}{2\rho}, \lambda(x)\right)$  and  $s \in [0, r]$ . We have*

1.  $\mu\left(\overline{\mathcal{B}}(x, r)\right) \geq c_9 r^d$ ,
2.  $\left| \frac{\mu(\overline{\mathcal{B}}(x, r))}{V_d r^d} - f(x) \right| \leq c_8 r$ ,
3.  $\mu\left(\overline{\mathcal{B}}(x, r) \setminus \overline{\mathcal{B}}(x, s)\right) \leq c_{10} r^{d-1} (r - s)$ .

with  $c_9 = f_{\min} J_{\min} V_d$ ,  $c_8 = c_7 + f_{\max} J_{\max} d 2^d \rho$  and  $c_{10} = d 2^d f_{\max} J_{\max} V_d$ .

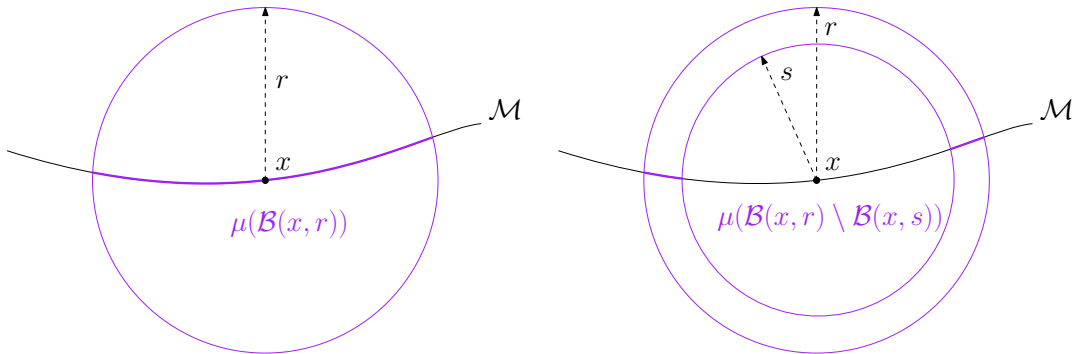


Figure 13: Representation of Proposition 3.13 Point 1 (left) and Point 3 (right).

*Proof.* Consider the map  $(\overline{\exp}_x^{\mathcal{M}})^{-1}$  defined in Equation (7) and the measure  $v_x = (\overline{\exp}_x^{\mathcal{M}})_*^{-1} \mu_x$  as defined in Lemma 3.11. In the following, we write  $T = T_x \mathcal{M}$ , and  $\overline{\mathcal{B}}_0 = (\overline{\exp}_x^{\mathcal{M}})^{-1}(\overline{\mathcal{B}}(x, r))$ .

*Point 1.* By definition of  $\mathbf{v}_x$ , we have  $\mu(\overline{\mathcal{B}}(x, r)) = \mathbf{v}_x(\overline{\mathcal{B}}_0)$ . Writing down the density  $g$  of  $\mathbf{v}_x$  yields

$$\mathbf{v}_x(\overline{\mathcal{B}}_0) = \int_{\overline{\mathcal{B}}_0} g(v) d\mathcal{H}^d(v).$$

According to the expression of  $g$  in Lemma 3.11, we have  $g \geq f_{\min} J_{\min}$ . Therefore,

$$\int_{\overline{\mathcal{B}}_0} g(v) d\mathcal{H}^d(v) \geq \int_{\overline{\mathcal{B}}_0} f_{\min} J_{\min} d\mathcal{H}^d(v) = f_{\min} J_{\min} \mathcal{H}^d(\overline{\mathcal{B}}_0).$$

Besides, since  $\overline{\mathcal{B}}_0 \supset \overline{\mathcal{B}}_T(0, r)$ , we have

$$\mathcal{H}^d(\overline{\mathcal{B}}_0) \geq \mathcal{H}^d(\overline{\mathcal{B}}_T(0, r)) = V_d r^d.$$

We finally obtain  $\mathbf{v}_x(\overline{\mathcal{B}}_0) \geq f_{\min} J_{\min} V_d r^d$ .

*Point 2.* Observe that  $\int_{\overline{\mathcal{B}}_T(0, r)} f(x) d\mathcal{H}^d(v) = f(x) V_d r^d$ . Hence

$$\begin{aligned} \left| \mu(\overline{\mathcal{B}}(x, r)) - f(x) V_d r^d \right| &= \left| \int_{\overline{\mathcal{B}}_0} g(v) d\mathcal{H}^d(v) - \int_{\overline{\mathcal{B}}_T(0, r)} f(x) d\mathcal{H}^d(v) \right| \\ &\leq \underbrace{\left| \int_{\overline{\mathcal{B}}_T(0, r)} (f(x) - g(v)) d\mathcal{H}^d(v) \right|}_A + \underbrace{\left| \int_{\overline{\mathcal{B}}_0 \setminus \overline{\mathcal{B}}_T(0, r)} g(v) d\mathcal{H}^d(v) \right|}_B. \end{aligned} \quad (8)$$

To bound Term A, notice that  $g(0) = f(\exp_x^{\mathcal{M}}(0)) J_0 = f(x)$ . Hence we can write:

$$\left| \int_{\overline{\mathcal{B}}_T(0, r)} (f(x) - g(v)) d\mathcal{H}^d(v) \right| \leq \int_{\overline{\mathcal{B}}_T(0, r)} |g(0) - g(v)| d\mathcal{H}^d(v).$$

Now, Lemma 3.11 gives  $|g(v) - g(0)| \leq c_7 r$ , and we eventually obtain the inequality  $\left| \int_{\overline{\mathcal{B}}_T(0, r)} (f(x) - g(v)) d\mathcal{H}^d(v) \right| \leq c_7 r V_d r^d$ .

On the other hand, we bound Term B thanks to the inclusion  $\overline{\mathcal{B}}_0 \subset \overline{\mathcal{B}}_T(0, c_5(\rho r))$ . Denote  $\mathcal{A} = \overline{\mathcal{B}}_T(0, c_5(\rho r)) \setminus \overline{\mathcal{B}}_T(0, r)$ . We have  $\overline{\mathcal{B}}_0 \setminus \overline{\mathcal{B}}_T(0, r) \subset \mathcal{A}$ , hence

$$\int_{\overline{\mathcal{B}}_0 \setminus \overline{\mathcal{B}}_T(0, r)} g(v) d\mathcal{H}^d(v) \leq \int_{\mathcal{A}} g(v) d\mathcal{H}^d(v) \leq f_{\max} J_{\max} \mathcal{H}^d(\mathcal{A}).$$

Moreover, we have

$$\mathcal{H}^d(\mathcal{A}) = \mathcal{H}^d(\overline{\mathcal{B}}_T(0, c_5(\rho r))) - \mathcal{H}^d(\overline{\mathcal{B}}_T(0, r)) = V_d (c_5(\rho r)^d - 1) r^d.$$

We can use  $c_5(\rho r) \leq 1 + 2\rho r \leq 2$  and the inequality  $a^d - 1 \leq d(a - 1)a^{d-1}$ , where  $a \geq 1$ , to get

$$(c_5(\rho r)^d - 1) \leq d \cdot (c_5(\rho r) - 1) \cdot c_5(\rho r)^{d-1} \leq d \cdot 2\rho r \cdot 2^{d-1}.$$

We finally deduce the following bound on Term B:

$$\int_{\overline{\mathcal{B}}_0 \setminus \overline{\mathcal{B}}_T(0, r)} g(v) d\mathcal{H}^d(v) \leq f_{\max} J_{\max} V_d r^d d \cdot \rho r 2^d.$$

Gathering Terms A and B, we obtain

$$\left| \mu(\overline{\mathcal{B}}(x, r)) - f(x) V_d r^d \right| \leq r (c_7 + f_{\max} J_{\max} d \rho 2^d) V_d r^d.$$

*Point 3.* Let us write

$$\begin{aligned} \mu(\overline{\mathcal{B}}(x, r) \setminus \overline{\mathcal{B}}(x, s)) &= \mathbf{v}_x((\exp_x^{\mathcal{M}})^{-1}(\overline{\mathcal{B}}(x, r) \setminus \overline{\mathcal{B}}(x, s))) \\ &= \int_{(\exp_x^{\mathcal{M}})^{-1}(\overline{\mathcal{B}}(x, r) \setminus \overline{\mathcal{B}}(x, s))} g(v) d\mathcal{H}^d(v). \end{aligned}$$

In spherical coordinates, this integral reads

$$\int_{(\exp_x \mathcal{M})^{-1}(\overline{\mathcal{B}}(x,r) \setminus \overline{\mathcal{B}}(x,s))} g(v) d\mathcal{H}^d(v) = \int_{v \in \partial \mathcal{B}_T(0,1)} \int_{t=a(v)}^{b(v)} g(tv) t^{d-1} dt dv, \quad (9)$$

where  $a(v)$  and  $b(v)$  are defined as follows: for every  $v \in T_x \mathcal{M}$  of unit norm, let  $\gamma_0$  be an arc-length parametrized geodesic with  $\gamma_0(0) = x$  and  $\dot{\gamma}_0(0) = v$ , and set  $a(v)$  and  $b(v)$  to be the first positive values such that  $\|\gamma(a(v)) - x\| = s$  and  $\|\gamma(b(v)) - x\| = r$ .

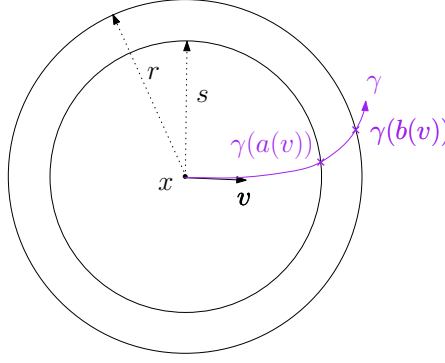


Figure 14: Illustration of  $a(v)$  and  $b(v)$  in Equation (9).

For any  $v \in \partial \mathcal{B}_T(0,1)$ , Lemma 2.4 Point 2 gives  $b(v) \leq 2r$ , hence

$$\int_{t=a(v)}^{b(v)} g(tv) t^{d-1} dt \leq \int_{t=a(v)}^{b(v)} f_{\max} J_{\max}(2r)^{d-1} dt.$$

Moreover, according to Lemma 2.4 Point 5, we have  $b(v) - a(v) \leq 2(r - s)$ , hence

$$\begin{aligned} \int_{t=a(v)}^{b(v)} f_{\max} J_{\max}(2r)^{d-1} dt &= (b(v) - a(v)) f_{\max} J_{\max}(2r)^{d-1} \\ &\leq 2(r - s) f_{\max} J_{\max}(2r)^{d-1}. \end{aligned}$$

From these last two equations we deduce

$$\begin{aligned} \int_{v \in \partial \mathcal{B}(0,1)} \int_{t=a(v)}^{b(v)} g(tv) t^{d-1} dt dv &\leq 2(r - s) f_{\max} J_{\max}(2r)^{d-1} \int_{v \in \partial \mathcal{B}(0,1)} dv \\ &= 2(r - s) f_{\max} J_{\max}(2r)^{d-1} dV_d. \end{aligned}$$

Going back to Equation (9), we obtain

$$\mu(\overline{\mathcal{B}}(x,r) \setminus \overline{\mathcal{B}}(x,s)) = 2^d dV_d f_{\max} J_{\max}(r-s) r^{d-1},$$

which concludes the proof.  $\square$

In Sect. 4, we will study the estimation of tangent spaces of  $\mathcal{M}$  thanks to the normal reach. By using the previous proposition, we will be able to give precise bounds around points  $x \in \mathcal{M}$  with large normal reach  $\lambda(x)$ . However, for points with small normal reach, we won't be able to use it. Therefore we need a version of Proposition 3.13 without normal reach condition. This is the aim of the following result.

**Proposition 3.14.** *Let  $x_0 \in \mathcal{M}_0$ ,  $r \leq \frac{1}{2\rho}$  and  $s \in [0, r]$ . We have*

1.  $\mu(\overline{\mathcal{B}}(x,r)) \geq c_9 r^d$
2.  $\mu(\overline{\mathcal{B}}(x,r) \setminus \overline{\mathcal{B}}(x,s)) \leq c_{11} r^{d-\frac{1}{2}} (r-s)^{\frac{1}{2}}$



with  $c_9 = f_{\min} J_{\min} V_d$  and  $c_{11} = \frac{f_{\max} J_{\max}}{f_{\min} J_{\min}} \left( \frac{\rho}{\sqrt{4-\sqrt{13}}} \right)^d d 2^{2d} \sqrt{3}$ .

Note that Point 1 is similar to Proposition 3.13 Point 1, and that Point 2 is a weaker form of Proposition 3.13 Point 3. There is no equivalent of Proposition 3.13 Point 2 without normal reach condition.

*Proof.* Let  $\mathcal{M}^x = \mathcal{M} \cap \overline{\mathcal{B}}(x, r)$  and  $\mathcal{M}_0^x = u^{-1}(\mathcal{M}^x)$ . Lemma 3.4 does not apply: it is not true that  $\mathcal{M}_0^x \subset \overline{\mathcal{B}}_{\mathcal{M}_0}(x_0, c_5(\rho r)r)$ . However, we can decompose  $\mathcal{M}_0^x$  in connected components  $C_0^i, i \in I$ . They are represented in Figure 15.

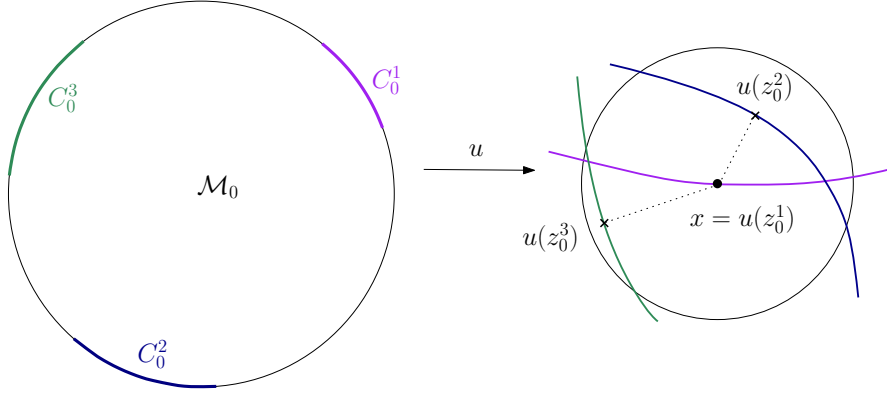


Figure 15: The connected components  $C_0^i, i \in I$ .

For every  $i \in I$ , let  $z_0^i$  be a minimizer of  $z_0 \mapsto \|z - x\|$  on  $C_0^i$ . We have  $x - z^i \perp T_{z^i} \mathcal{M}$ . Following the same proof as Lemma 3.4, one shows that  $C_0^i$  is included in the geodesic ball  $\overline{\mathcal{B}}_{\mathcal{M}_0}(z_0^i, \frac{1}{\rho})$ . Hence we can consider  $\mu_0^i$ , the measure  $\mu_0$  restricted to  $C_0^i$ , and define  $\nu_0^i = (\overline{\exp}_{z_0^i})_*^{-1} \mu_0^i$ , as in Remark 3.12. The measure  $\nu_0^i$  admits  $g_0^i$  as a density over the  $d$ -dimensional Hausdorff measure on  $T_{z_0^i} \mathcal{M}_0$ , where

$$g_0^i(\nu) = f_0(\overline{\exp}_{z_0^i} \mathcal{M}_0(\nu)) \cdot J_\nu \cdot 1_{(\overline{\exp}_{z_0^i})^{-1}(C_0^i)}(\nu).$$

*Point 1.* By definition of  $\mu$ , can write

$$\mu(\overline{\mathcal{B}}(x, r)) = \mu_0(u^{-1}(\overline{\mathcal{B}}(x, r))) = \sum_{i \in I} \mu_0(C_0^i).$$

Denote by  $0 \in I$  the index of the connected component of  $\mathcal{M}_0^x$  which contains  $x_0$ . We have  $C_0^0 \supset \overline{\mathcal{B}}_{\mathcal{M}_0}(x_0, r)$ . As in the proof of Proposition 3.13 Point 1, we deduce that

$$\begin{aligned} \mu_0(C_0^0) &\geq \int_{(\overline{\exp}_{z_0^0})^{-1}(C_0^0)} g_0^0 \cdot d\mathcal{H}^d \\ &\geq f_{\min} J_{\min} \mathcal{H}^d((\overline{\exp}_{z_0^0})^{-1}(C_0^0)) = f_{\min} J_{\min} V_d r^d. \end{aligned}$$

Therefore,  $\mu(\overline{\mathcal{B}}(x, r)) \geq f_{\min} J_{\min} V_d r^d$ .

*Point 2.* For any  $i \in I$ , define  $D_0^i = C_0^i \cap u^{-1}(\overline{\mathcal{B}}(x, r) \setminus \overline{\mathcal{B}}(x, s))$ . Let us show that

$$\mu_0(D_0^i) \leq f_{\max} J_{\max} 2^{d-1} \sqrt{6} V_d \cdot r^{d-1} \sqrt{r^2 - s^2}. \quad (10)$$

As in Equation (9), we write this measure as

$$\int_{(\overline{\exp}_{z_0^i})^{-1}(D_0^i)} g_0^i(y) d\mathcal{H}^d(y) = \int_{v \in \partial \mathcal{B}(0,1)} \int_{t=a(v)}^{b(v)} g_0^i(tv) t^{d-1} dt dv,$$

where  $a(v)$  and  $b(v)$  are defined as in the proof of Proposition 3.13 Point 3: for every  $v \in \subset T_{z^i} \mathcal{M}$  of unit norm, let  $\gamma_0$  be an arc-length parametrized geodesic with  $\gamma_0(0) = z_0^i$  and  $\dot{\gamma}_0(0) = v$ , and set  $a(v)$  and  $b(v)$  to be the first positive values such that  $\|\gamma(a(v)) - x\| \geq s$  and  $\|\gamma(b(v)) - x\| = r$ . For any  $v \in \partial \mathcal{B}_T(0, 1)$ , Lemma 2.4 Point 2 gives  $b(v) \leq 2r$ , and Lemma 2.4 Point 4 gives  $b(v) - a(v) \leq \sqrt{6}\sqrt{r^2 - s^2}$ . We deduce that

$$\int_{t=a(v)}^{b(v)} f_{\max} J_{\max}(2r)^{d-1} dt \leq \sqrt{6}\sqrt{r^2 - s^2} f_{\max} J_{\max}(2r)^{d-1}.$$

Therefore,

$$\int_{v \in \partial \mathcal{B}(0,1)} \int_{t=a(v)}^{b(v)} g_0^i(tv) t^{d-1} dt dv \leq \sqrt{6}\sqrt{r^2 - s^2} f_{\max} J_{\max}(2r)^{d-1} dV_d,$$

which yields Equation (10).

We now gather the connected components  $D_0^i$ . Since  $u^{-1}(\overline{\mathcal{B}}(x, r) \setminus \overline{\mathcal{B}}(x, s)) = \bigcup_{i \in I} D_0^i$ , we have

$$\mu(\overline{\mathcal{B}}(x, r) \setminus \overline{\mathcal{B}}(x, s)) = \sum_{i \in I} \mu_0(D_i).$$

Using Equation (10) we get

$$\mu(\overline{\mathcal{B}}(x, r) \setminus \overline{\mathcal{B}}(x, s)) \leq |I| f_{\max} J_{\max} 2^{d-1} \sqrt{6} dV_d \cdot r^{d-1} \sqrt{r^2 - s^2},$$

where  $|I|$  is the cardinal of  $I$ . Let us show that  $|I| \leq \frac{1}{f_{\min} J_{\min} V_d} \left(\frac{2\rho}{\alpha}\right)^d$ , with  $\alpha = \sqrt{4 - \sqrt{13}}$ , which will conclude the proof.

Let  $i, j \in I$  such that  $i \neq j$ . We first show that  $d_{\mathcal{M}_0}(z_0^i, z_0^j) \geq \frac{\alpha}{\rho}$ . Let  $\gamma_0: [0, T] \rightarrow \mathcal{M}_0$  be a geodesic from  $z_0^i$  to  $z_0^j$ . Consider the map  $\phi: t \mapsto \|\gamma(t) - x\|$ . Since  $C_0^i$  and  $C_0^j$  are disjoint connected components, there must be a  $t^* < T$  such that  $\|\gamma(t^*) - x_0\| > r$ . Moreover, according to Lemma 2.4 Point 1,  $\phi$  is increasing on  $[0, T_1]$  where  $T_1 = \frac{\sqrt{2}}{\rho} \sqrt{2 - \sqrt{3 + \rho^2 l^2}}$ . Since  $\phi(T) \leq r$ , we deduce that  $T$  is greater than  $T_1$ . Note that the assumption  $r \leq \frac{1}{2\rho}$  yields  $T_2 \geq \frac{\alpha}{\rho}$ . Hence we obtain the bound

$$d_{\mathcal{M}_0}(z_0^i, z_0^j) = T \geq T_1 \geq \frac{\alpha}{\rho}.$$

This implies that the geodesic balls  $\mathcal{B}_{\mathcal{M}_0}\left(z_0^i, \frac{\alpha}{2\rho}\right)$ ,  $i \in I$ , are disjoint. Therefore,

$$1 \geq \mu_0\left(\bigcup_i \mathcal{B}_{\mathcal{M}_0}\left(z_0^i, \frac{\alpha}{2\rho}\right)\right) \geq |I| f_{\min} J_{\min} V_d \left(\frac{\alpha}{2\rho}\right)^d,$$

and we deduce that  $|I| \leq \frac{1}{f_{\min} J_{\min} V_d} \left(\frac{2\rho}{\alpha}\right)^d$ . □

### 3.3 Sublevel sets of the normal reach

In this subsection, we assume the Hypotheses 2 and 3, as well as Hypothesis 1', stated in the next paragraph. This last hypothesis can be seen as a strengthening of Hypothesis 1. Our goal is to give an upper bound on  $\mu_0(\lambda_0^t)$ , the measure of the set of points  $x_0 \in \mathcal{M}_0$  with normal reach not greater than  $t$  (see Proposition 3.19). This proves a result announced in Subsect. 2.4: Hypothesis 4 is a consequence of Hypotheses 1', 2 and 3. We close this subsection with a remark concerning generalizations of this result. Since Hypothesis 4 trivially holds when the immersion is an embedding (with  $r_4 = \min \lambda_0$  and  $c_4 = 0$ ), we shall also suppose that  $u$  is not an embedding.

First, we say that a finite collection  $A$  of linear subspaces of  $E$  is in *general position* if

$$\text{codim}\left(\bigcap_{V \in A} V\right) = \sum_{V \in A} \text{codim}(V),$$

where we define  $\text{codim}(V) = \dim(E) - \dim(V)$ . Now, we say that the immersion  $u: \mathcal{M}_0 \rightarrow \mathcal{M}$  is *self-transverse* (also called *completely regular* in [51]) if for any point  $x \in \mathcal{M}$ , the collection of tangent spaces  $\{T_y \mathcal{M} \mid y_0 \in \mathcal{M}_0, x = y\}$  is in *general position*. Suppose that  $u$  is self-transverse, and denote by  $\mathcal{N}_0$  be the self-intersections of  $\mathcal{M}_0$ :

$$\mathcal{N}_0 = \{x_0 \in \mathcal{M}_0 \mid \exists y_0 \in \mathcal{M}_0, x_0 \neq y_0, x = y\}.$$

Its image is denoted  $\mathcal{N} = u(\mathcal{N}_0)$ . Equivalently,  $\mathcal{N}_0$  is the set of points with zero normal reach, that is,  $\mathcal{N}_0 = \lambda_0^{-1}(\{0\})$ . We also have  $\mathcal{N} = \lambda^{-1}(\{0\})$ . In general,  $\mathcal{N}_0$  and  $\mathcal{N}$  are not submanifolds, but only (closed) immersed manifolds. The subset  $\mathcal{N}_0$  can be decomposed as a disjoint union

$$\mathcal{N}_0 = \bigsqcup_{i \geq 2} \mathcal{N}_0^{(i)} \quad \text{where} \quad \mathcal{N}_0^{(i)} = \{x_0 \in \mathcal{M}_0 \mid |u^{-1}(\{x\})| = i\},$$

and where  $|\cdot|$  denotes the cardinal. In other words,  $\mathcal{N}_0^{(i)}$  is the set of points of  $\mathcal{M}_0$  whose image is shared by exactly  $i$  distinct points of  $\mathcal{M}_0$ . Each  $\mathcal{N}_0^{(i)}$  is a submanifold of  $\mathcal{M}_0$ , not necessarily closed, of dimension  $i \dim(\mathcal{M}_0) - (i-1) \dim(E)$  [51, Lemma 2.3]. Moreover, the tangent spaces of  $\mathcal{N}^{(i)} = u(\mathcal{N}_0^{(i)})$  can be described as:

$$T_x \mathcal{N}^{(i)} = \bigcap_{y_0 \in u^{-1}(x)} T_{y_0} \mathcal{M}. \quad (11)$$

In order to state the proofs of this subsection, we shall make the following assumption: the immersion  $u$  only has *double points*, that is,  $\mathcal{N}_0$  is equal to  $\mathcal{N}_0^{(2)}$ . We shall refer to this assumption as

**Hypothesis 1'.** The immersion  $u$  is self-transverse, and only has double points.

In this case,  $\mathcal{N}_0$  is a submanifold of  $\mathcal{M}_0$ , of dimension  $2 \dim(\mathcal{M}_0) - \dim(E)$ . Most of the examples we will consider later in the paper satisfy this assumption. They are curves in the plane (Examples 5.1, 5.4 and 5.6) or surfaces in the space (Examples 5.2, 5.7).

We will also need a few quantities related to the immersion. Let  $\mathcal{D}_0$  be the set of critical points of the Euclidean distance on  $\mathcal{M}_0$ , that is,

$$\mathcal{D}_0 = \{(x_0, y_0) \in \mathcal{M}_0 \times \mathcal{M}_0 \mid x_0 \neq y_0, x - y \perp T_{y_0} \mathcal{M} \text{ and } x - y \perp T_{x_0} \mathcal{M}\}. \quad (12)$$

Also, let  $\mathcal{C}_0$  be the set of double points of  $\mathcal{M}_0$ :

$$\mathcal{C}_0 = \{(x_0, y_0) \in \mathcal{M}_0 \times \mathcal{M}_0 \mid x_0 \neq y_0 \text{ and } x = y\}. \quad (13)$$

Note that the projection of  $\mathcal{C}_0$  on the first coordinate is  $\mathcal{N}_0$ . Moreover, we have  $\mathcal{C}_0 \subset \mathcal{D}_0$ , and these sets are compact. Since  $\mathcal{C}_0$  is an isolated subset of  $\mathcal{D}_0$  by Lemma 2.3 Point 4, we have that  $\mathcal{D}_0 \setminus \mathcal{C}_0$  also is compact. Consider the quantity

$$\Delta = \inf \{\|x - y\| \mid (x_0, y_0) \in \mathcal{D}_0 \setminus \mathcal{C}_0\}. \quad (14)$$

The constant  $\Delta$  can be understood as the minimal length of the nonzero bottlenecks of  $\mathcal{M}$ . From the compactness of  $\mathcal{D}_0 \setminus \mathcal{C}_0$  we deduce that  $\Delta > 0$ . Moreover, we define

$$\Delta_0 = \inf \{\|x - y\| \mid x_0 \in \mathcal{N}_0, y_0 \in \mathcal{M}_0, x \neq y, x - y \perp T_{y_0} \mathcal{M}\}. \quad (15)$$

It is a measure of regularity around the self-intersections of  $\mathcal{M}$ . Using Lemma 2.3 Point 5, one proves that this infimum is taken over a compact set, hence that  $\Delta_0 > 0$ . Last, we will need a measure a similarity between linear subspaces. If  $U, V$  denote two linear subspaces of  $E$ , let their minimal angle be

$$\angle(U, V) = \inf \left\{ \arccos \left( \frac{|\langle u, v \rangle|}{\|u\| \|v\|} \right) \mid u \in U, v \in V, u, v \in (U \cap V)^\perp \right\},$$

where  $\inf \emptyset = 0$  by convention. Note that  $\angle(U, V) > 0$  when  $U \neq V$ . Now, define

$$\Theta = \inf \{ \angle(T_x \mathcal{M}, T_y \mathcal{M}) \mid (x_0, y_0) \in \mathcal{C}_0 \}. \quad (16)$$

According to the self-transversality hypothesis and the compactness of  $\mathcal{C}_0$ , we have  $\Theta > 0$ . These constants are represented in Figure 16.

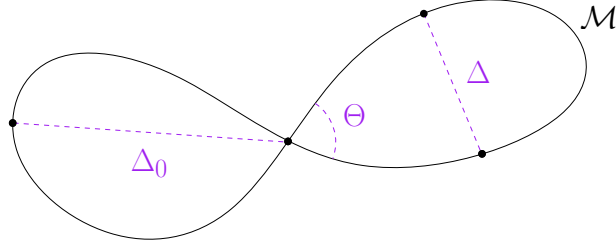


Figure 16: The constants  $\Delta$ ,  $\Delta_0$  and  $\Theta$  associated to  $\mathcal{M}$ .

In order to bound the measure  $\mu_0(\lambda_0^t)$ , we will prove that the sublevel set  $\lambda_0^t$  is included in a thickening of  $\mathcal{N}_0$ . By bounding the measure of this thickening, we will obtain the main result (Proposition 3.19). We start with a lemma which describes the situation around self-intersection points of  $\mathcal{M}_0$ .

**Lemma 3.15.** *Let  $(x_0, y_0) \in \mathcal{C}_0$  (defined in Equation (13)). Let  $\gamma_0: I \rightarrow \mathcal{M}_0$  (resp.  $\gamma'_0$ ) be an arc-length parametrized geodesic starting from  $x_0$  (resp. from  $y_0$ ), and denote  $v = \dot{\gamma}(0)$  (resp.  $v' = \dot{\gamma}'(0)$ ). Let  $\theta = \arccos(|\langle v, v' \rangle|)$  be their angle. Let  $\delta, \delta' \geq 0$  such that  $\delta' \leq \delta \leq \frac{\sin(\theta)}{2\rho}$ . Then we have*

$$\|\gamma(\delta) - \gamma'(\delta')\| \geq \frac{\sin(\theta)}{2} \delta.$$

As a consequence, if  $v$  is orthogonal to  $T_x \mathcal{M} \cap T_y \mathcal{M}$ , then the distance from  $\gamma(\delta)$  to  $u(\overline{\mathcal{B}}_{\mathcal{M}_0}(y_0, \delta))$  is lower bounded by  $\frac{\sin(\Theta)}{2} \delta$ .

*Proof.* Let us introduce  $\bar{x} = x + \delta v$  and  $\bar{y} = y + \delta' v'$ , as represented in Figure 17.

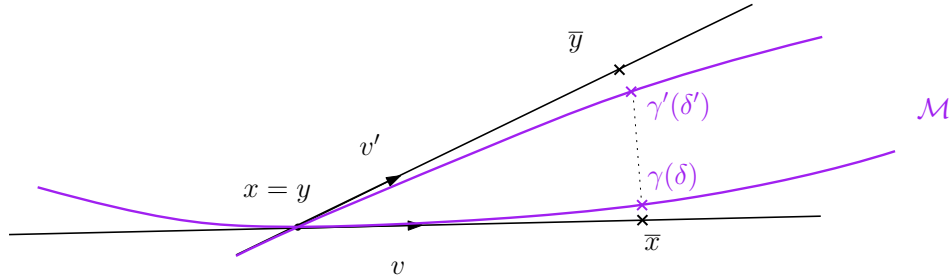


Figure 17: Situation in Lemma 3.15.

The triangle inequality yields

$$\|\gamma(\delta) - \gamma'(\delta')\| \geq \|\bar{x} - \bar{y}\| - \|\gamma(\delta) - \bar{x}\| - \|\gamma'(\delta') - \bar{y}\|.$$

According to Lemma 2.3 Point 1, we have  $\|\gamma(\delta) - \bar{x}\| \leq \frac{\rho}{2} \delta^2$  and  $\|\gamma'(\delta') - \bar{y}\| \leq \frac{\rho}{2} \delta'^2 \leq \frac{\rho}{2} \delta^2$ . Moreover,  $\|\bar{x} - \bar{y}\|$  is not lower than  $\|\bar{x} - z\|$ , where  $z$  is the projection of  $\bar{x}$  on the line spanned by  $v'$ . Elementary trigonometry shows that  $\|\bar{x} - z\| = \sin(\theta) \delta$ . Hence the previous equation yields

$$\|\gamma(\delta) - \gamma'(\delta')\| \geq \sin(\theta) \delta - \frac{\rho}{2} \delta^2 - \frac{\rho}{2} \delta^2 = \sin(\theta) \delta \left( 1 - \frac{\rho}{\sin(\theta)} \delta \right),$$

and we conclude with  $\delta \leq \frac{\sin(\theta)}{2\rho}$ . □

The following lemma shows that, around  $\mathcal{N}$ , the immersed manifold  $\mathcal{M}$  is a union of two transversally intersecting pieces.

**Lemma 3.16.** *For any  $r < \min\left(\frac{1}{2\rho}, \Delta_0\right)$  and  $x \in \mathcal{N}$ , the set  $u^{-1}(\overline{\mathcal{B}}(x, r))$  is made up of two connected components, and we have*

$$u^{-1}(\overline{\mathcal{B}}(x, r)) \subset \bigcup_{y_0 \in u^{-1}(\{x\})} \overline{\mathcal{B}}_{\mathcal{M}_0}(y_0, 2r).$$

*Proof.* Consider  $\mathcal{M}_0^x = u^{-1}(\overline{\mathcal{B}}(x, r))$  and  $C_0^i$ ,  $i \in I$ , its connected components, as represented in Figure 18. Let us denote  $C_0^0$  the connected component that contains  $x_0$ .

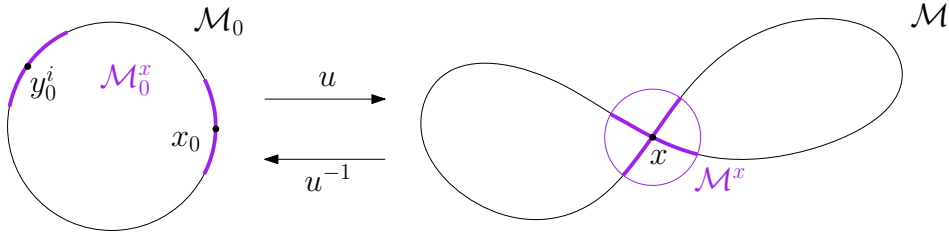


Figure 18: Situation in Lemma 3.16.

For any  $i \in I \setminus \{0\}$ , let  $y_0^i$  be a minimizer of  $y_0 \mapsto \|y - x\|$  on  $C_0^i$ . It satisfies  $x - y^i \perp T_{y^i} \mathcal{M}$ . Since  $r$  has been chosen lower than  $\Delta_0$  (defined in Equation (15)), we must have  $y^i = x$ , that is to say,  $y_0^i \in u^{-1}(\{x\})$ . Using that  $u^{-1}(\{x\})$  consists of two elements, we deduce that  $\mathcal{M}_0^x$  is made up of two connected components.

Now, as we have seen in the proof of Lemma 3.4, these connected components satisfy  $C_0^i \subset \overline{\mathcal{B}}_{\mathcal{M}_0}(y_0^i, c_5(\rho r))$ . The result follows from  $c_5(\rho r) < 2$ .  $\square$

We can now connect the normal reach to the distance to  $\mathcal{N}_0$ . We first prove that, close to a self-intersection point, the normal reach is lower bounded by the geodesic distance to that point.

**Lemma 3.17.** *Let  $x_0 \in \mathcal{M}_0$  and denote by  $\delta = d_{\mathcal{M}_0}(x_0, \mathcal{N}_0)$  the geodesic distance from  $x_0$  to  $\mathcal{N}_0$ . Suppose that  $\delta < \min\left(\frac{1}{4\rho}, \frac{\sin(\Theta)}{2}, \frac{\Delta_0}{2}\right)$ . Then  $\lambda_0(x_0) \geq \frac{\sin(\Theta)}{2} \delta$ .*

*Proof.* Let  $y_0$  be a projection of  $x_0$  on  $\mathcal{N}_0$ , that is, a point that minimizes the geodesic distance  $d_{\mathcal{M}_0}(x_0, y_0)$  with  $y_0 \in \mathcal{N}_0$ . Let  $\gamma_0: I \rightarrow \mathcal{M}_0$  be a geodesic from  $y_0$  to  $x_0$ , and denote  $v_0 = \dot{\gamma}_0(0)$ . Note that  $v \perp T_{y_0} \mathcal{N}$ , otherwise  $y_0$  would not be a minimizer. Denote  $r = \|x - y\|$  and  $\mathcal{M}_0^y = u^{-1}(\overline{\mathcal{B}}(y, 2r))$ . Let also  $y'_0$  be the other point of  $\mathcal{M}_0$  such that  $y' = y$ .

First, let us show that  $\lambda_0(x_0) \leq r$ . According to Lemma 3.16,  $\mathcal{M}_0^y$  consists of two connected components,  $C_0$  that contains  $y_0$ , and  $C'_0$  that contains  $y'_0$ . Now, consider a minimizer of  $z_0 \mapsto \|z - x\|$  on  $C'_0$ . This point satisfies  $x - z \perp T_z \mathcal{M}$  and  $x_0 \neq z_0$ . Thus,

$$\lambda_0(x_0) \leq \|x - z\| \leq \|x - y\| \leq r,$$

as announced. Moreover, using the inequality  $\overline{\mathcal{B}}(x, r) \subset \overline{\mathcal{B}}(y, 2r)$ , we deduce that  $z$  realizes the normal reach of  $x_0$ , that is,  $\lambda_0(x_0) = \|x - z\|$ .

To conclude, consider the geodesic  $\gamma_0$  defined above. We have seen that  $v$  is orthogonal to the tangent space  $T_{y_0} \mathcal{N}$ . Since  $T_{y_0} \mathcal{N} = T_{x_0} \mathcal{M} \cap T_{y_0} \mathcal{M}$  by Equation (11), we have  $v \perp (T_{x_0} \mathcal{M} \cap T_{y_0} \mathcal{M})$ , hence we can apply the consequence of Lemma 3.15 to get

$$\lambda_0(x_0) = \|x - z\| = \|\gamma(\delta) - z\| \geq \frac{\sin(\Theta)}{2} \delta,$$

as wanted.  $\square$

The following lemma is a converse of Lemma 3.17: points with small normal reach are close to the self-intersection submanifold  $\mathcal{N}_0$ .

**Lemma 3.18.** *Let  $x_0 \in \mathcal{M}_0$  such that  $\lambda_0(x_0) \leq \min\left(\frac{\sin(\Theta)}{8\rho}, \frac{\sin(\Theta)^2}{4}, \frac{\Delta_0 \sin(\Theta)}{4}, \Delta\right)$ . Then  $\lambda_0(x_0) \geq \frac{\sin(\Theta)}{2} d_{\mathcal{M}_0}(x_0, \mathcal{N}_0)$ .*

*Proof.* Put  $r = \lambda_0(x_0)$ , and denote the sublevel set  $\lambda_0^r = \lambda_0^{-1}([0, r])$ . Let  $C_0$  denote the connected component of  $x_0$  in  $\lambda_0^r$ . We have seen in Remark 3.2 that the normal reach  $\lambda_0$  is lower semi-continuous. Hence  $C_0$  is closed, and  $\lambda_0$  attains a minimum on it. Let  $y_0$  be a minimizer of  $\lambda_0$  on  $C_0$ . Let us prove that  $\lambda_0(y_0) = 0$  by contradiction. Suppose that  $\lambda_0(y_0) > 0$ , and let  $z_0 \in \mathcal{M}_0$  be such that  $z_0 \neq y_0$ ,  $y - z \perp T_{z_0} \mathcal{M}$  and  $\lambda_0(y_0) = \|y - z\|$ . These points are represented in the following figure.

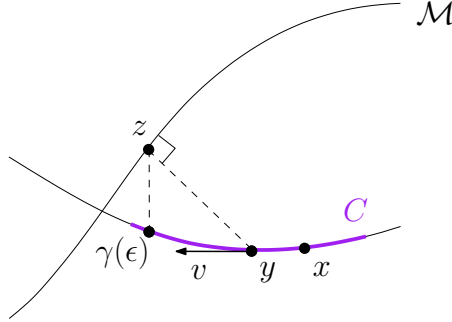


Figure 19: Situation in Lemma 3.18, supposing by contradiction that  $\lambda_0(y_0) > 0$ .

Since  $\|y - z\| = \lambda_0(y_0) \leq \lambda_0(x_0) < \Delta$ , where  $\Delta$  has been defined in Equation (14), the vector  $y - z$  is not orthogonal to  $T_y \mathcal{M}$ . Let  $v_0 \in T_{y_0} \mathcal{M}_0$  such that  $\langle v, z - y \rangle > 0$ , and consider an arc-length parametrized geodesic  $\gamma: I \rightarrow \mathcal{M}_0$  with  $\gamma_0(0) = y_0$  and  $\dot{\gamma}_0(0) = v_0$ . We will show that  $\lambda_0(\gamma_0(\epsilon)) < \lambda_0(y_0)$  for  $\epsilon > 0$  small enough, contradicting the minimality of  $y_0$ .

On the one hand, for  $\epsilon > 0$  small enough, the bound  $\langle v, z - y \rangle > 0$  yields

$$\|z - \gamma(\epsilon)\| < \|z - y\| = \lambda_0(y_0) \quad (17)$$

On the other hand, Lemma 2.3 Point 5 and  $y - z \perp T_{z_0} \mathcal{M}$  gives  $d_{\mathcal{M}_0}(z_0, y_0) \geq \frac{1}{\rho}$ . Together with the assumption  $\lambda_0(x_0) < \frac{\sin(\Theta)}{8\rho} < \frac{1}{4\rho}$ , we deduce that  $d_{\mathcal{M}_0}(z_0, y_0) > 4\lambda_0(x_0)$ . By continuity, for  $\epsilon > 0$  small enough, we also have

$$d_{\mathcal{M}_0}(z_0, \gamma_0(\epsilon)) > 4\lambda_0(x_0). \quad (18)$$

Now, we deduce from  $\lambda_0(x_0) \geq \lambda_0(y_0)$  and Equation (17) and (18) that  $d_{\mathcal{M}_0}(z_0, \gamma_0(\epsilon)) > 4\|z - \gamma(\epsilon)\|$ . Therefore we can apply Lemma 3.5 on  $z_0$  and  $\gamma_0(\epsilon)$  to get

$$\lambda_0(\gamma_0(\epsilon)) \leq \|z - \gamma(\epsilon)\| < \lambda_0(y_0),$$

which contradicts the minimality of  $y_0$ . We conclude that  $\lambda_0(y_0) = 0$ .

Next, let  $\gamma_0: [0, T] \rightarrow \mathcal{M}_0$  be a path from  $y_0$  to  $x_0$  in  $C_0$ . Let us show that, for all  $t \in [0, T]$ ,

$$\lambda_0(\gamma_0(t)) \geq \frac{\sin(\Theta)}{2} d_{\mathcal{M}_0}(\gamma_0(t), \mathcal{N}_0). \quad (19)$$

According to Lemma 3.17, it is enough to show that  $d_{\mathcal{M}_0}(\gamma_0(t), \mathcal{N}_0) < c$ , where  $c = \min\left(\frac{1}{4\rho}, \frac{\sin(\Theta)}{2}, \frac{\Delta_0}{2}\right)$ . By contradiction, suppose that  $d_{\mathcal{M}_0}(\gamma_0(t), \mathcal{N}_0) \geq c$  for some  $t$ . Since  $d_{\mathcal{M}_0}(\gamma_0(0), \mathcal{N}_0) = d_{\mathcal{M}_0}(y_0, \mathcal{N}_0) = 0$ ,

we can consider the first value  $t \in [0, T]$  such that  $d_{\mathcal{M}_0}(\gamma_0(t), \mathcal{N}_0) = c$ . Lemma 3.17 then gives  $\lambda_0(\gamma_0(t)) \geq \frac{\sin(\Theta)}{2}c$ . Besides, by definition of  $C_0$ , we have  $\lambda_0(\gamma_0(t)) \leq \lambda_0(x_0)$ . We deduce that

$$\lambda_0(x_0) \geq \frac{\sin(\Theta)}{2} \cdot \min\left(\frac{1}{4\rho}, \frac{\sin(\Theta)}{2}, \frac{\Delta_0}{2}\right),$$

which contradicts the assumptions of the lemma. We now obtain the result from Equation (19) at  $t = T$ .  $\square$

We now prove the main result of this subsection.

**Proposition 3.19.** *Suppose that the immersion  $u$  satisfies Hypotheses 1', 2 and 3. Let  $\alpha = \dim(E) - \dim(\mathcal{M}_0)$ . For every  $r < r_{13}$ , we have*

$$\mu_0(\lambda_0^r) \leq c_{13}r^\alpha + O(r^{\alpha+1}),$$

where  $r_{13} = \min\left(\frac{\sin(\Theta)}{8\rho}, \frac{\sin(\Theta)^2}{4}, \frac{\Delta_0 \sin(\Theta)}{4}, \Delta\right)$  and  $c_{13} = \left(\frac{2}{\sin(\Theta)}\right)^\alpha V_\alpha f_{\max} \mathcal{H}_{\mathcal{M}_0}^{d'}(\mathcal{N}_0)$ .

*Proof.* Let  $d'$  be the dimension of  $\mathcal{N}_0$ , and  $\alpha$  its codimension in  $\mathcal{M}_0$ . Since  $\dim(\mathcal{N}_0) = 2\dim(\mathcal{M}_0) - \dim(E)$  by Hypothesis 1', we have  $\alpha = \dim(E) - \dim(\mathcal{M}_0)$ . Besides, according to Lemma 3.18, the sublevel set  $\lambda_0^r$  is included in the geodesic thickening  $\mathcal{N}_0^t$  of  $\mathcal{N}_0$ , where  $t = \frac{2}{\sin(\Theta)}r$ . According to Weyl's Tube Formula [52, Theorem 9.23], the volume of  $\mathcal{N}_0^t$  is

$$\mathcal{H}_{\mathcal{M}_0}^{d'}(\mathcal{N}_0^t) = V_\alpha \cdot \mathcal{H}_{\mathcal{M}_0}^{d'}(\mathcal{N}_0) \cdot t^\alpha + O(t^{\alpha+1}).$$

where  $\mathcal{H}_{\mathcal{M}_0}^{d'}(\mathcal{N}_0)$  is the  $d'$ -dimensional volume of  $\mathcal{N}_0$ , and  $V_\alpha$  the volume of the unit ball in  $\mathbb{R}^\alpha$ . Using the density of  $\mu_0$  given by Hypothesis 2, we can write  $\mu_0(\mathcal{N}_0^t) \leq f_{\max} \mathcal{H}_{\mathcal{M}_0}^{d'}(\mathcal{N}_0^t)$ . Hence

$$\mu_0(\mathcal{N}_0^t) \leq f_{\max} \mathcal{H}_{\mathcal{M}_0}^{d'}(\mathcal{N}_0^t) = f_{\max} \mathcal{H}_{\mathcal{M}_0}^{d'}\left(\mathcal{N}_0^{\frac{2}{\sin(\Theta)}r}\right),$$

and the result follows.  $\square$

*Remark 3.20.* It seems reasonable to think that Proposition 3.19 is still valid when replacing Hypothesis 1' with the weaker Hypothesis 1, at least with  $\alpha = 1$ . In this case, the quantities  $\Delta$ ,  $\Delta_0$  and  $\Theta$  are still well-defined, and Lemmas 3.16 and 3.15 hold. However, Lemmas 3.17 and 3.18 may not be true anymore. An illustration of this is given by the two following intersecting surfaces of  $\mathbb{R}^4$ :

$$\mathcal{M}^{(1)} = \{(a, b, 0, 0) \mid a, b \in \mathbb{R}\} \quad \text{and} \quad \mathcal{M}^{(2)} = \{(a, 0, c, a^2) \mid a, c \in \mathbb{R}\}.$$

These submanifolds intersect at  $\mathcal{M}^{(1)} \cap \mathcal{M}^{(2)} = \{0\}$ , but their intersection is not transverse, since their tangent spaces does not span the fourth canonical basis vector of  $\mathbb{R}^4$ . The distance from a point  $x = (a, 0, 0, a^2)$  of  $\mathcal{M}^{(2)}$  to  $\mathcal{M}^{(1)}$  is  $a^2$ . Besides, for small values of  $a$ , the geodesic distance from  $x$  to 0 on  $\mathcal{M}^{(2)}$  is approximately  $a$ . But there is no constant  $c$  such that  $a^2 \geq c \cdot a$ . Hence Lemma 3.17 does not hold anymore. Nonetheless, the subset of point of  $\mathcal{M}^{(2)}$  at distance at most  $r$  from  $\mathcal{M}^{(1)}$  is

$$\left\{(a, 0, c, a^2) \in \mathcal{M}^{(2)} \mid a^4 + c^2 \leq r^2\right\}.$$

This is approximately the rectangle  $\{(a, 0, c, a^2) \in \mathcal{M}^{(2)} \mid |a| \leq \sqrt{r}, |c| \leq r\}$ , whose volume is  $r^{\frac{3}{2}}$ . We see here that Proposition 3.19 holds with  $\alpha = \frac{3}{2}$ .

*Remark 3.21.* We can also see that Proposition 3.19 is true under the following conditions: Hypotheses 1, 2, 3 and  $\dim(\mathcal{M}_0) = 1$ . Indeed, in this case, the self-intersecting manifold  $\mathcal{N}_0$  is a finite subset of  $\mathcal{M}_0$ . Knowing this fact, the proofs of Lemmas 3.15, 3.16, 3.17, 3.18 and Proposition 3.19 can be used without modification. In this case, the result reads  $\mu_0(\lambda_0^r) \leq c_{13}r + O(r^2)$ .

## 4 Tangent space estimation

We now come back to our original setting:  $\mathcal{M}_0$  is a manifold of dimension  $d \geq 1$ , immersed in  $E = \mathbb{R}^n$  via  $u: \mathcal{M}_0 \rightarrow \mathcal{M}$ . Moreover,  $\mathcal{M}_0$  is endowed with a measure  $\mu_0$ . The push-forward measure is denoted  $\mu = u_*\mu_0$ . In this section, we show that one can estimate the tangent spaces of  $\mathcal{M}$ , based on the measure  $\mu$ , or a close measure  $\nu$ , via the computation of local covariance matrices.

### 4.1 Local covariance matrices and lifted measures

We remind the reader that the aim of this work is to estimate the homotopy type of  $\mathcal{M}_0$ , or its homology groups, from the measure  $\nu$ . As explained in the introduction and in Subsect. 2.4, our strategy consists in estimating the *lifted manifold*

$$\check{\mathcal{M}} = \left\{ \left( x, \frac{1}{d+2} p_{T_x \mathcal{M}} \right) \mid x_0 \in \mathcal{M}_0 \right\},$$

where  $p_{T_x \mathcal{M}}$  is the matrix of the orthogonal projection on the tangent space  $T_x \mathcal{M}$ , seen as an element of  $M(E)$ , the space of  $n \times n$  matrices. The normalization term  $\frac{1}{d+2}$  has been chosen in accordance with Proposition 4.1, stated in the next subsection, and makes our method independant of the dimension  $d$ . Note that the set  $\check{\mathcal{M}}$  can also be described as the image of  $\mathcal{M}_0$  under the map

$$\check{u}: x_0 \mapsto \left( x, \frac{1}{d+2} p_{T_x \mathcal{M}} \right).$$

Using Hypothesis 1, we deduce that  $\check{\mathcal{M}}$  is diffeomorphic to  $\mathcal{M}_0$ , hence that their homotopy types and homology groups coincide. Adopting a measure theoretical point of view on the problem, we will actually estimate the *exact lifted measure*, defined as the push-forward  $\check{\mu}_0 = \check{u}_*\mu_0$ . We explain in Subsect. 5.1 how one can infer the homotopy type and homology groups of  $\mathcal{M}_0$  from  $\check{\mu}_0$ . Note that this measure can also be defined as

$$\check{\mu}_0 = (u_*\mu_0)(x_0) \otimes \left\{ \delta_{\frac{1}{d+2} p_{T_x \mathcal{M}}} \right\} \quad (20)$$

by disintegration of measure. Here is another alternative definition of  $\check{\mu}_0$ : for any smooth  $\phi: E \times M(E) \rightarrow \mathbb{R}$  with compact support,

$$\int \phi(x, A) d\check{\mu}_0(x, A) = \int \phi \left( u(x_0), \frac{1}{d+2} p_{T_x \mathcal{M}} \right) d\mu_0(x_0). \quad (21)$$

In order to approximate  $\check{\mu}_0$ , we have to propose an estimator of the tangent spaces. We consider the following construction. If  $x$  in any vector of  $E = \mathbb{R}^n$ , seen as a row vector, the *tensor product* is defined as the  $n \times n$  matrix  $x^{\otimes 2} = x^t \cdot x$ .

**Definition 4.1.** Let  $\nu$  be any probability measure on  $E$ . Let  $r > 0$  and  $x \in \text{supp}(\nu)$ . The *local covariance matrix of  $\nu$  around  $x$  at scale  $r$*  is the following matrix:

$$\Sigma_\nu(x) = \int_{\mathcal{B}(x, r)} (x - y)^{\otimes 2} \frac{d\nu(y)}{\nu(\mathcal{B}(x, r))}.$$

We also define the *normalized local covariance matrix* as  $\bar{\Sigma}_\nu(x) = \frac{1}{r^2} \Sigma_\nu(x)$ .

Note that  $\Sigma_\nu(x)$  and  $\bar{\Sigma}_\nu(x)$  depend on  $r$ , which is not made explicit in the notation. The normalization factor  $\frac{1}{r^2}$  of the normalized local covariance matrix is justified by Proposition 4.1. Moreover, we introduce the following notations: for every  $r > 0$  and  $x \in \text{supp}(\nu)$ ,

- $\nu_x$  is the restriction of  $\nu$  to the ball  $\mathcal{B}(x, r)$ ,
- $\bar{\nu}_x = \frac{1}{\nu(\mathcal{B}(x, r))} \nu_x$  is the corresponding probability measure.



Thus the local covariance matrix can be written as  $\Sigma_v(x) = \int (x-y)^{\otimes 2} d\bar{v}_x(y)$ .

We note that such notions have already been studied in the context of Topological Data Analysis. The collection of probability measures  $\{\bar{v}_x\}_{x \in \text{supp}(v)}$  is called in [39, Sect. 3.3] the local truncation of  $v$  at scale  $r$ . The map  $x \mapsto \Sigma_v(x)$  is called in [38, Sect. 2.2] the multiscale covariance tensor field of  $v$  associated to the truncation kernel.

We now propose an estimator of the lifted measure  $\check{\mu}_0$ , inspired by Equations (20) and (21).

**Definition 4.2.** For any measure  $v$  on  $E$ , we denote by  $\check{v}$  the measure on  $E \times M(E)$  defined by

$$\check{v} = v(x) \otimes \left\{ \delta_{\bar{\Sigma}_v(x)} \right\}.$$

It is called the *lifted measure* associated to  $v$ . In other words, for every smooth  $\phi: E \times M(E) \rightarrow \mathbb{R}$  with compact support, we have

$$\int \phi(x, A) d\check{v}(x, A) = \int \phi\left(x, \bar{\Sigma}_v(x)\right) dv(x).$$

In accordance with the local covariance matrices, the lifted measure  $\check{v}$  depends on the parameter  $r$  which is not made explicit in the notation.

In order to compare these measures, we consider a Wasserstein-type distance on the space  $E \times M(E)$ . Fix  $\gamma > 0$ , and let  $\|\cdot\|_\gamma$  be the Euclidean norm on  $E \times M(E)$  defined as

$$\|(x, A)\|_\gamma^2 = \|x\|^2 + \gamma^2 \|A\|_F^2, \quad (22)$$

where  $\|\cdot\|$  represents the usual Euclidean norm on  $E$  and  $\|\cdot\|_F$  represents the Frobenius norm on  $M(E)$ . Let  $p \geq 1$ . We denote by  $W_{p,\gamma}(\cdot, \cdot)$  the  $p$ -Wasserstein distance with respect to this metric. By definition, if  $\alpha, \beta$  are probability measures on  $E \times M(E)$ , then  $W_{p,\gamma}(\alpha, \beta)$  can be written as

$$W_{p,\gamma}(\alpha, \beta) = \inf_{\pi} \left( \int_{(E \times M(E))^2} \|(x, A) - (y, B)\|_\gamma^p d\pi((x, A), (y, B)) \right)^{\frac{1}{p}}, \quad (23)$$

where the infimum is taken over all measures  $\pi$  on  $(E \times M(E))^2$  with marginals  $\alpha$  and  $\beta$ .

The parameter  $\gamma$  of the norm  $\|\cdot\|_\gamma$  has been designed to balance the importance given to the Euclidean information ( $E$ -coordinate) and matrix information ( $M(E)$ -coordinate) in  $E \times M(E)$ . The more  $\gamma$  is large, the more the matrix information will be relatively important. Since there is no canonical choice of  $\gamma$ , it will remain as a free parameter in the rest of the paper. In the experiments of the next section, we will choose the value  $\gamma = 1$  or  $2$ , for it seemed relevant in practice. For a discussion about how this parameter may influence the persistent homology of the lifted manifold  $\check{\mathcal{M}}$ , we refer the reader to [53, Subsect. 4.4].

We subdivide the rest of this section in four subsections. They respectively consists in showing that

- **Consistency:** if  $\mu_0$  is a measure satisfying the Hypotheses 2 and 3, then  $W_{p,\gamma}(\check{\mu}_0, \check{\mu})$  is small (Proposition 4.2),
- **Stability of the localized measures:** in addition, if  $v$  is a measure on  $E$  such that  $W_p(\mu, v)$  is small, then so is  $W_1(\bar{\mu}_y, \bar{v}_y)$  (Lemmas 4.7 and 4.8),
- **Stability of the lifted measures:** consequently,  $W_{p,\gamma}(\check{\mu}, \check{v})$  is also small (Proposition 4.11)
- **Approximation:** under the previous hypotheses,  $W_{p,\gamma}(\check{\mu}_0, \check{v})$  is small (Theorem 4.14).

These measures fit in a commutative diagram:

$$\begin{array}{ccccc} \mathcal{M}_0 & \xrightarrow{\check{u}} & E \times M(E) & & \\ & \searrow u & \swarrow \text{proj} & & \\ & & E & & \end{array} \quad \begin{array}{ccccc} \mu_0 & \xrightarrow{\check{u}_*} & \check{\mu}_0 & \xrightarrow{\check{v}} & \check{v} \\ & \searrow u_* & \nearrow g_* & \nearrow (f_v)_* & \\ & & \mu & \xrightarrow{(f_\mu)_*} & v \end{array}$$

where the maps  $g, f_\mu$  and  $f_\nu: E \rightarrow E \times \mathbf{M}(E)$  are defined as

$$g: x \mapsto \left(x, \frac{1}{d+2} p_{T_x \mathcal{M}}\right), \quad f_\mu: x \mapsto \left(x, \bar{\Sigma}_\mu(x)\right), \quad f_\nu: x \mapsto \left(x, \bar{\Sigma}_\nu(x)\right).$$

Note that the map  $g$  is well-defined only on points  $x \in \mathcal{M}$  that are not self-intersection points, i.e., points  $x$  such that  $\lambda(x) > 0$ . Under Hypothesis 4,  $g$  is well-defined  $\mu$ -almost surely. The maps  $f_\mu$  and  $f_\nu$  are defined respectively on  $\text{supp}(\mu)$  and  $\text{supp}(\nu)$ .

## 4.2 Consistency of the estimation

In this subsection, we assume that  $\mathcal{M}_0$  and  $\mu_0$  satisfy Hypotheses 2 and 3. We first show that the normalized covariance matrix approximates the tangent spaces of  $\mathcal{M}$ , as long as the parameter  $r$  is chosen smaller than the normal reach. A similar result appears in [29, Lemma 13] in the case where  $\mathcal{M}$  is a submanifold and  $\mu$  is the uniform distribution on  $\mathcal{M}$ . Based on this result, we deduce that the lifted measure  $\check{\mu}$  is close to the exact lifted measure  $\check{\mu}_0$ . The quality of this approximation depends on the measure of the set of points with small normal reach, i.e., points where the tangent spaces are not well-estimated.

**Proposition 4.1.** *Let  $x_0 \in \mathcal{M}_0$  and  $r < \min\left(\lambda(x), \frac{1}{2\rho}\right)$ . Denote by  $p_{T_x \mathcal{M}}$  the orthogonal projection matrix on the tangent space  $T_x \mathcal{M}$ . We have*

$$\left\| \bar{\Sigma}_\mu(x) - \frac{1}{d+2} p_{T_x \mathcal{M}} \right\|_F \leq c_{14} r,$$

where  $c_{14} = 6\rho + 4 \frac{c_7}{f_{\min} J_{\min}} + \frac{f_{\max}}{f_{\min} J_{\min}} 2^d d \rho + \frac{c_8}{f_{\min} J_{\min}}$ .

*Proof.* According to [29, Lemma 11], the matrix  $r^2 \frac{1}{d+2} p_{T_x \mathcal{M}}$  is equal to

$$\Sigma_* = \int_{\bar{\mathcal{B}}_{T_x \mathcal{M}}(0, r)} y^{\otimes 2} \cdot \frac{d\mathcal{H}^d(y)}{V_d r^d}.$$

Hence the proposition reduces to  $\|\Sigma_\mu(x) - \Sigma_*\|_F \leq c_{14} r^3$ . Let us write  $T = T_x \mathcal{M}$ ,  $\bar{\mathcal{B}} = \bar{\mathcal{B}}(x, r)$  and  $\bar{\mathcal{B}}_0 = (\bar{\text{exp}}_x^\mathcal{M})^{-1}(\bar{\mathcal{B}})$ , where  $(\bar{\text{exp}}_x^\mathcal{M})^{-1}$  as been defined in Equation (7). We consider the following intermediate matrices:

$$\begin{aligned} \Sigma_1 &= \int_{\bar{\mathcal{B}}} \left( (\bar{\text{exp}}_x^\mathcal{M})^{-1}(x') \right)^{\otimes 2} d\bar{\mu}_x(x'), \\ \Sigma_2 &= \int_{\bar{\mathcal{B}}_0} g(0) y^{\otimes 2} \cdot \frac{d\mathcal{H}^d(y)}{|\mu_x|}, \\ \Sigma_3 &= \int_{\bar{\mathcal{B}}_T(0, r)} g(0) y^{\otimes 2} \cdot \frac{d\mathcal{H}^d(y)}{|\mu_x|}. \end{aligned}$$

The triangle inequality now yields:

$$\|\Sigma_\mu(x) - \Sigma_*\|_F \leq \underbrace{\|\Sigma_\mu(x) - \Sigma_1\|_F}_A + \underbrace{\|\Sigma_1 - \Sigma_2\|_F}_B + \underbrace{\|\Sigma_2 - \Sigma_3\|_F}_C + \underbrace{\|\Sigma_3 - \Sigma_*\|_F}_D. \quad (24)$$

*Term A.* By definition of the local covariance matrix, we have

$$\Sigma_\mu(x) = \int_{\bar{\mathcal{B}}(x, r)} (x - x')^{\otimes 2} \bar{\mu}_x(x').$$

We use the upper bound

$$\begin{aligned}\|\Sigma_\mu(x) - \Sigma_1\|_F &\leq \int_{\overline{\mathcal{B}}(x,r)} \left\| (x-x')^{\otimes 2} - \left( (\overline{\text{exp}}_x^\mathcal{M})^{-1}(x') \right)^{\otimes 2} \right\|_F d\overline{\mu}_x(x') \\ &\leq \sup_{x' \in \mathcal{M} \cap \overline{\mathcal{B}}(x,r)} \left\| (x-x')^{\otimes 2} - \left( (\overline{\text{exp}}_x^\mathcal{M})^{-1}(x') \right)^{\otimes 2} \right\|_F.\end{aligned}$$

Let  $x' \in \mathcal{M} \cap \overline{\mathcal{B}}(x,r)$ . According to Lemma 3.4, we have  $\left\| (\overline{\text{exp}}_x^\mathcal{M})^{-1}(x') \right\| \leq 2r$ . Moreover,  $\|x - x'\| \leq r$ , and Lemma 4.3, stated in the following subsection, gives

$$\left\| (x-x')^{\otimes 2} - \left( (\overline{\text{exp}}_x^\mathcal{M})^{-1}(x') \right)^{\otimes 2} \right\|_F \leq (r+2r) \left\| (x'-x) - \left( (\overline{\text{exp}}_x^\mathcal{M})^{-1}(x') \right) \right\|. \quad (25)$$

Now, let us justify that

$$\left\| (x'-x) - \left( (\overline{\text{exp}}_x^\mathcal{M})^{-1}(x') \right) \right\| \leq \frac{\rho}{2} d_{\mathcal{M}_0}(x_0, x'_0)^2. \quad (26)$$

If we write  $x' = \gamma(\delta)$  with  $\gamma$  a geodesic such that  $\gamma(0) = x$  and  $\delta = d_{\mathcal{M}_0}(x_0, x'_0)$ , then  $(\overline{\text{exp}}_x^\mathcal{M})^{-1}(x') = \delta \dot{\gamma}(0)$ , and we get

$$\left\| (x'-x) - \left( (\overline{\text{exp}}_x^\mathcal{M})^{-1}(x') \right) \right\| = \|\gamma(\delta) - (x + \delta \dot{\gamma}(0))\| \leq \frac{\rho}{2} \delta^2,$$

where we used Lemma 2.3 Point 1 for the last inequality. Hence Equation (26) is true. Combined with Lemma 3.4, which gives  $d_{\mathcal{M}_0}(x_0, x'_0) \leq 2\|x - x'\| \leq 2r$ , we obtain

$$\left\| (x-x')^{\otimes 2} - \left( (\overline{\text{exp}}_x^\mathcal{M})^{-1}(x') \right)^{\otimes 2} \right\|_F \leq \frac{\rho}{2} (2r)^2 = 2\rho r^2.$$

We now use Equation (25) to deduce  $\|\Sigma_\mu(x) - \Sigma_1\|_F \leq (r+2r)2\rho r^2 = 6\rho r^3$ .

*Term B.* By transfer, we can write  $\Sigma_1$  as

$$\Sigma_1 = \int_{\overline{\mathcal{B}}} \left( (\overline{\text{exp}}_x^\mathcal{M})^{-1}(x') \right)^{\otimes 2} \frac{d\mathcal{H}^d(y)}{|\mu_x|} = \int_{\overline{\mathcal{B}_0}} g(y) y^{\otimes 2} \cdot \frac{d\mathcal{H}^d(y)}{|\mu_x|}.$$

We deduce the upper bound

$$\|\Sigma_1 - \Sigma_2\|_F \leq \int_{\overline{\mathcal{B}_0}} |g(0) - g(y)| \|y^{\otimes 2}\| \frac{d\mathcal{H}^d(y)}{|\mu_x|}.$$

According to Lemma 4.3,  $\|y^{\otimes 2}\| = \|y\|^2 \leq (2r)^2$ , and Lemma 3.11 gives  $|g(y) - g(0)| \leq c_7 r$ . Therefore,

$$\|\Sigma_1 - \Sigma_2\|_F \leq 4r^2 \cdot c_7 r \cdot \frac{\mathcal{H}^d(\overline{\mathcal{B}_0})}{|\mu_x|}.$$

To conclude, note that  $|\mu_x| \geq f_{\min} J_{\min} \mathcal{H}^d(\overline{\mathcal{B}_0})$  by Proposition 3.13 Point 1, hence we obtain  $\|\Sigma_1 - \Sigma_2\|_F \leq 4 \frac{c_7}{f_{\min} J_{\min}} r^3$ .

*Term C.* As for the previous terms, we use the upper bound

$$\|\Sigma_2 - \Sigma_3\|_F \leq \int_{\overline{\mathcal{B}_T(0,r)} \setminus \overline{\mathcal{B}_0}} \|g(0) \cdot y^{\otimes 2}\|_F \frac{d\mathcal{H}^d(y)}{|\mu_x|}.$$

On the one hand,  $\|g(0) \cdot y^{\otimes 2}\|_F \leq g(0) \cdot r^2 \leq f_{\max} r^2$ , and we get

$$\|\Sigma_2 - \Sigma_3\|_F \leq f_{\max} r^2 \frac{\mathcal{H}^d(\overline{\mathcal{B}}_T(0, r) \setminus \overline{\mathcal{B}}_0)}{|\mu_x|}.$$

On the other hand, since  $\overline{\mathcal{B}}_0 \subset \overline{\mathcal{B}}_T(x, c_5(\rho r)r)$ , we have

$$\mathcal{H}^d(\overline{\mathcal{B}}_0 \setminus \overline{\mathcal{B}}_T(0, r)) = (c_5(\rho r)r)^d V_d - r^d V_d.$$

The inequality  $a^d - 1 \leq d(a - 1)a^{d-1}$ , where  $a \geq 1$ , gives

$$(c_5(\rho r)r)^d V_d - r^d V_d \leq V_d r^d \cdot d(c_5(\rho r) - 1)2^{d-1}.$$

Combined with the inequalities  $c_5(\rho r) \leq 1 + 2\rho r$  and  $|\mu_x| \geq f_{\min} J_{\min} V_d r^d$ , we get

$$\|\Sigma_2 - \Sigma_3\|_F \leq \frac{f_{\max}}{f_{\min} J_{\min}} 2^d d \rho r^3.$$

*Term D.* Let us write  $\Sigma^*$  as

$$\Sigma_* = \int_{\overline{\mathcal{B}}_{T_x, \mathcal{M}}(0, r)} y^{\otimes 2} \cdot \frac{|\mu_x|}{V_d r^d} \cdot \frac{d\mathcal{H}^d(y)}{|\mu_x|}.$$

Hence we have

$$\|\Sigma_3 - \Sigma_*\|_F \leq \int_{\overline{\mathcal{B}}_T(0, r)} \left| \frac{|\mu_x|}{V_d r^d} - f(x) \right| \|y^{\otimes 2}\|_F \frac{d\mathcal{H}^d(y)}{|\mu_x|}.$$

According to Proposition 3.13 Point 2,  $\left| \frac{|\mu_x|}{V_d r^d} - f(x) \right| \leq c_8 r$ . Moreover,  $\|y^{\otimes 2}\|_F \leq r^2$  and  $\int_{\overline{\mathcal{B}}_T(0, r)} \frac{d\mathcal{H}^d(y)}{|\mu_x|} \leq \frac{1}{f_{\min} J_{\min}}$ . Therefore,  $\|\Sigma_3 - \Sigma_*\|_F \leq \frac{c_8}{f_{\min} J_{\min}} r^3$ . We deduce the result by summing Terms A, B, C and D.  $\square$

We now deduce a result concerning the lifted measures  $\check{\mu}$  and  $\check{\mu}_0$  (defined in Subsect. 4.1). We remind the reader that the notation  $\lambda^r$  refers to the sublevel set  $\lambda^{-1}([0, r])$ . Hence the quantity  $\mu(\lambda^r)$  is the measure of the set of points  $x \in \mathcal{M}$  such that  $\lambda(x) \leq r$ .

**Proposition 4.2.** *Let  $r < \frac{1}{2p}$ . Then*

$$W_{p, \gamma}(\check{\mu}, \check{\mu}_0) \leq \gamma \left( 2\mu(\lambda^r)^{\frac{1}{p}} + c_{14} r \right).$$

*Proof.* Define the map  $\phi: \mathcal{M}_0 \rightarrow (E \times \mathbf{M}(E)) \times (E \times \mathbf{M}(E))$  as

$$\phi: x_0 \mapsto \left( \left( x, \bar{\Sigma}_\mu(x) \right), \left( x, \frac{1}{d+2} p_{T_x, \mathcal{M}} \right) \right),$$

and consider the measure  $\pi = \phi_* \mu_0$ . It is a transport plan between  $\check{\mu}$  and  $\check{\mu}_0$ . By definition of the Wasserstein distance,

$$W_{p, \gamma}^p(\check{\mu}, \check{\mu}_0) \leq \int \| (x, T) - (x', T') \|_\gamma^p d\pi((x, T), (x', T')),$$

hence we can use this transport plan and write

$$\begin{aligned} W_{p, \gamma}^p(\check{\mu}, \check{\mu}_0) &\leq \int \left\| \left( x, \frac{1}{r^2} \Sigma_\mu(x) \right) - \left( x, \frac{1}{d+2} p_{T_x, \mathcal{M}} \right) \right\|_\gamma^p d\mu(x) \\ &= \gamma^p \int \left\| \frac{1}{r^2} \Sigma_\mu(x) - \frac{1}{d+2} p_{T_x, \mathcal{M}} \right\|_F^p d\mu(x). \end{aligned}$$

We split this last integral into the sets  $A = \lambda^r$  and  $B = E \setminus \lambda^r$ .

On  $A$ , we use the upper bound  $\left\| \frac{1}{r^2} \Sigma_\mu(x) - \frac{1}{d+2} p_{T_x \mathcal{M}} \right\|_F \leq \left\| \frac{1}{r^2} \Sigma_\mu(x) \right\|_F + \left\| \frac{1}{d+2} p_{T_x \mathcal{M}} \right\|_F \leq 1 + 1$  to obtain

$$\int_A \left\| \frac{1}{r^2} \Sigma_\mu(x) - \frac{1}{d+2} p_{T_x \mathcal{M}} \right\|_F^p d\mu(x) \leq 2^p \mu(A).$$

On  $B$ , we use Proposition 4.1 to get

$$\int_B \left\| \frac{1}{r^2} \Sigma_\mu(x) - \frac{1}{d+2} p_{T_x \mathcal{M}} \right\|_F^p d\mu(x) \leq (c_{14}r)^p.$$

Combining these two inequalities yields  $W_{p,\gamma}^p(\check{\mu}, \check{\mu}_0) \leq \gamma^p (2^p \mu(A) + (c_{14}r)^p)$ . Using the inequality  $(a+b)^{\frac{1}{p}} \leq a^{\frac{1}{p}} + b^{\frac{1}{p}}$ , where  $a, b \geq 0$ , we deduce

$$W_{p,\gamma}(\check{\mu}, \check{\mu}_0) \leq \gamma \left( 2\mu(A)^{\frac{1}{p}} + c_{14}r \right),$$

which is the result.  $\square$

### 4.3 Stability of localization of measures

This technical subsection is dedicated to proving stability results for localization of measures. Throughout the subsection, we consider two measures  $\mu$  and  $\nu$  on  $E$ . We show that, under some hypotheses on  $\mu$ , an upper bound on the Wasserstein distance  $W_p(\mu, \nu)$  gives an upper bound on the the Wassertsein distance  $W_1(\bar{\mu}_y, \bar{\nu}_y)$  between their localized measures (see Lemmas 4.7 and 4.8). We close this subsection with a comment about the sharpness of our bounds.

The results of this subsection only rely on the following hypotheses about  $\mu$ :

**Hypothesis 5.**  $\exists c_9 > 0, \forall x \in \text{supp}(\mu), \forall t \in [0, \frac{1}{2\rho})$ ,

$$\mu(\bar{\mathcal{B}}(x, t)) \geq c_9 t^d.$$

**Hypothesis 6.**  $\exists c_{10} > 0, \forall x \in \text{supp}(\mu), \exists \lambda(x) \geq 0, \forall s, t \in [0, \min(\lambda(x), \frac{1}{2\rho})]$  s.t.  $s \leq t$ ,

$$\mu(\bar{\mathcal{B}}(x, t) \setminus \bar{\mathcal{B}}(x, s)) \leq c_{10} t^{d-1} (t-s).$$

**Hypothesis 7.**  $\exists c_{11} > 0, \forall x \in \text{supp}(\mu), \forall s, t \in [0, \frac{1}{2\rho})$  s.t.  $s \leq t$ ,

$$\mu(\bar{\mathcal{B}}(x, t) \setminus \bar{\mathcal{B}}(x, s)) \leq c_{11} t^{d-\frac{1}{2}} (t-s)^{\frac{1}{2}}.$$

Note that these Hypotheses 5, 6 and 7 are consequences of the initial Hypotheses 2 and 3. Indeed, as stated in Propositions 3.13 and 3.14, these new hypotheses hold with  $\lambda(x)$  being the normal reach of  $\mathcal{M}$  at  $x$ , and with the constants  $c_9 = f_{\min} J_{\min} V_d$ ,  $c_{10} = d 2^d f_{\max} J_{\max} V_d$  and

$$c_{11} = \frac{f_{\max} J_{\max}}{f_{\min} J_{\min}} \left( \frac{\rho}{\sqrt{4 - \sqrt{13}}} \right)^d d 2^{2d} \sqrt{3}.$$

In order to state the results of this subsection in a more general setting, we will only invoke the Hypotheses 5, 6 and 7. We first state a lemma that will be useful in what follows.

**Lemma 4.3.** *For every  $x, y \in E$ , we have  $\|x^{\otimes 2} - y^{\otimes 2}\|_F \leq (\|x\| + \|y\|) \|x - y\|$ .*

*Proof.* We apply the triangle inequality to  $x^t x - y^t y = (x - y)^t x + y^t (x - y)$ :

$$\begin{aligned} \|x^t x - y^t y\|_F &\leq \|(x - y)^t x\|_F + \|y^t (x - y)\|_F \leq \|x - y\| \|x\| + \|y\| \|x - y\| \\ &= (\|x\| + \|y\|) \|x - y\|, \end{aligned}$$

which gives the bound.  $\square$

Next, let us compare a measure and its submeasures. If  $\mu$  is a measure of positive mass (potentially with  $|\mu| \neq 1$ ), we remind the reader that the notation  $\bar{\mu}$  refers to the corresponding probability measure  $\frac{1}{\mu(E)}\mu$ . Moreover, a *submeasure* of  $\mu$  is a measure  $\mu'$  such that for all measurable set  $A \subset E$ , we have  $\mu'(A) \leq \mu(A)$ .

**Lemma 4.4.** *Let  $\mu$  be any measure of positive mass, and let  $\mu'$  be a submeasure of  $\mu$  with  $|\mu'| > 0$ . Suppose that  $\text{supp}(\mu)$  is included in a ball  $\mathcal{B}(x, r)$ . Then*

$$W_p(\bar{\mu}, \bar{\mu}') \leq 2 \left(1 - \frac{|\mu'|}{|\mu|}\right)^{\frac{1}{p}} r.$$

*In particular, if  $\mu$  is a probability measure, then  $W_p(\mu, \bar{\mu}') \leq 2(1 - |\mu'|)^{\frac{1}{p}} r$ .*

*Proof.* We start with the second inequality. Consider the intermediate probability measure  $\omega = \mu' + (1 - |\mu'|)\delta_x$ , where  $\delta_x$  is the Dirac mass (represented in Figure 20). We shall use the triangle inequality  $W_p(\mu, \mu') \leq W_p(\mu, \omega) + W_p(\omega, \mu')$ . We can write

- $\mu = \mu' + (\mu - \mu')$ ,
- $\omega = \mu' + (1 - |\mu'|)\delta_x$ ,
- $\bar{\mu}' = \mu' + (\bar{\mu}' - \mu')$ .

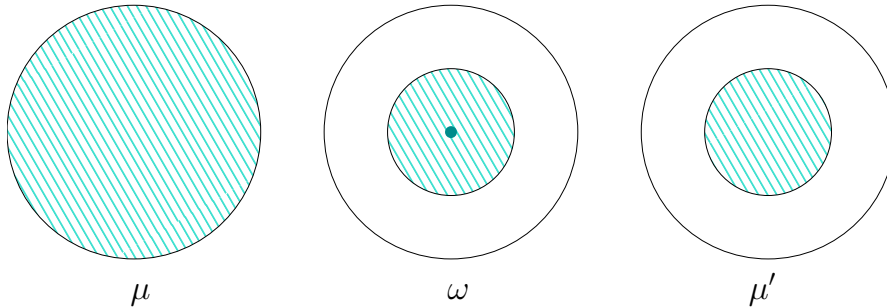


Figure 20: The measures involved in the proof of Lemma 4.4. A hatched area represents the support of the measure, and a point represents a Dirac mass.

Observe that  $\mu$  and  $\omega$  admits  $\mu'$  as a common submeasure of mass  $|\mu'|$ . Therefore we can build a transport plan between  $\mu$  and  $\omega$  where only a mass  $1 - |\mu'|$  of  $\mu$  is moved to  $x$ . In other words,

$$W_p(\mu, \omega) \leq (1 - |\mu'|)^{\frac{1}{p}} r.$$

Similarly, one shows that  $W_p(\omega, \bar{\mu}') \leq (1 - |\mu'|)^{\frac{1}{p}} r$ .

Now let us prove the first inequality. Since  $\mu'$  is a submeasure of  $\mu$  of mass  $|\mu'|$ , then  $\frac{1}{|\mu|}\mu'$  is a submeasure of  $\bar{\mu} = \frac{1}{|\mu|}\mu$  of mass  $\frac{1}{|\mu|}|\mu'|$ . We then apply the previous inequality.  $\square$

We now compare the localized measures of  $\mu$ , (defined in Subsect. 4.1).

**Lemma 4.5.** Let  $x \in \text{supp}(\mu)$ . Suppose that  $x$  satisfies Hypotheses 5 and 6 with  $r < \min\left(\lambda(x), \frac{1}{2p}\right)$ . Let  $y \in E$  such that  $\|x - y\| < \frac{r}{4}$ . Then  $|\mu_x| > 0$ ,  $|\mu_y| > 0$  and

$$W_1(\overline{\mu_x}, \overline{\mu_y}) \leq c_{15} \|x - y\|,$$

with  $c_{15} = 2 \left(1 + 4 \frac{5^{d-1}}{3^d}\right) \frac{c_{10}}{c_9}$ .

*Proof.* It is clear that  $|\mu_y| > 0$  since  $\mu(\overline{\mathcal{B}}(y, r)) \geq \mu(\overline{\mathcal{B}}(x, r - \|x - y\|))$  and  $x \in \text{supp}(\mu)$ . Let us show the inequality  $W_1(\overline{\mu_x}, \overline{\mu_y}) \leq c_{15} \|x - y\|$  by studying the measure  $\mu$  on the intersection  $\overline{\mathcal{B}}(x, r) \cap \overline{\mathcal{B}}(y, r)$ . Let  $\mu_{x,y}$  be the restriction of  $\mu$  to  $\overline{\mathcal{B}}(x, r) \cap \overline{\mathcal{B}}(y, r)$ , and  $\overline{\mu_{x,y}}$  the corresponding probability measure. The triangle inequality gives:

$$W_1(\overline{\mu_x}, \overline{\mu_y}) \leq \underbrace{W_1(\overline{\mu_x}, \overline{\mu_{x,y}})}_A + \underbrace{W_1(\overline{\mu_{x,y}}, \overline{\mu_y})}_B. \quad (27)$$

*Term A.* Let us show that  $W_1(\overline{\mu_x}, \overline{\mu_{x,y}}) \leq 2 \frac{c_{10}}{c_9} \|x - y\|$ . Note that  $\overline{\mu_{x,y}}$  is a submeasure of  $\overline{\mu_x}$ . According to Lemma 4.4, we have

$$W_1(\overline{\mu_x}, \overline{\mu_{x,y}}) \leq 2 \left(1 - \frac{|\mu_{x,y}|}{|\mu_x|}\right) r = 2 \frac{|\mu_x| - |\mu_{x,y}|}{|\mu_x|} r.$$

We know from Hypothesis 5 that  $|\mu_x| \geq c_9 r^d$ . On the other hand,

$$\begin{aligned} |\mu_x| - |\mu_{x,y}| &= \mu(\overline{\mathcal{B}}(x, r)) - \mu(\overline{\mathcal{B}}(x, r) \cap \overline{\mathcal{B}}(y, r)) \\ &\leq \mu(\overline{\mathcal{B}}(x, r)) - \mu(\overline{\mathcal{B}}(x, r - \|x - y\|)), \end{aligned}$$

hence we can apply Hypothesis 6 to get  $|\mu_x| - |\mu_{x,y}| \leq c_{10} r^{d-1} \|x - y\|$ . We finally obtain

$$W_1(\overline{\mu_x}, \overline{\mu_{x,y}}) \leq 2 \frac{c_{10} r^{d-1} \|x - y\|}{c_9 r^d} r = 2 \frac{c_{10}}{c_9} \|x - y\|.$$

*Term B.* Similarly, Lemma 4.4 yields

$$W_1(\overline{\mu_y}, \overline{\mu_{x,y}}) \leq 2 \frac{|\mu_y| - |\mu_{x,y}|}{|\mu_y|} r.$$

Let us show that we still have  $|\mu_y| \geq a' r^d$  and  $|\mu_y| - |\mu_{x,y}| \leq b' r^{d-1} \|x - y\|$  with the constants  $a' = (\frac{3}{4})^d c_9$  and  $b' = 2(\frac{5}{4})^{d-1} c_{10}$ . The first inequality comes from Hypothesis 5:

$$\mu(\overline{\mathcal{B}}(y, r)) \geq \mu(\overline{\mathcal{B}}(x, r - \|x - y\|)) \geq c_9 (r - \|x - y\|)^d$$

and  $\|x - y\| \leq \frac{r}{4}$ . The second inequality comes from Hypothesis 6:

$$\begin{aligned} \mu(\overline{\mathcal{B}}(y, r)) - \mu(\overline{\mathcal{B}}(x, r) \cap \overline{\mathcal{B}}(y, r)) &\leq \mu(\overline{\mathcal{B}}(x, r + \|x - y\|)) - \mu(\overline{\mathcal{B}}(x, r - \|x - y\|)) \\ &\leq c_{10} (r + \|x - y\|)^{d-1} 2 \|x - y\| \end{aligned}$$

and  $\|x - y\| \leq \frac{r}{4}$ . To conclude,

$$W_1(\overline{\mu_y}, \overline{\mu_{x,y}}) \leq 2 \frac{2(\frac{5}{4})^{d-1} r^{d-1} c_9 \|x - y\|}{2(\frac{3}{4})^d c_{10} r^d} r = 8 \frac{5^{d-1}}{3^d} \frac{c_{10}}{c_9} \|x - y\|.$$

We obtain the result by summing Terms A and B.  $\square$

The following lemma is the counterpart of Lemma 4.5 when replacing Hypothesis 6 with the weaker Hypothesis 7.

**Lemma 4.6.** Let  $x \in \text{supp}(\mu)$ . Suppose that  $x$  satisfies Hypotheses 5 and 7 at  $x$  with  $r < \frac{1}{2\rho}$ . Let  $y \in E$  such that  $\|x - y\| < \frac{r}{4}$ . Then  $|\mu_x|, |\mu_y| > 0$ , and

$$W_1(\overline{\mu_x}, \overline{\mu_y}) \leq c_{16} r^{\frac{1}{2}} \|x - y\|^{\frac{1}{2}},$$

$$\text{with } c_{16} = \left(2 + \frac{2^{\frac{5}{2}} 5^{d-\frac{1}{2}}}{3^d}\right) \frac{c_{11}}{c_9}.$$

*Proof.* The proof is similar to Lemma 4.5 with slight modifications. We still consider

$$W_1(\overline{\mu_x}, \overline{\mu_y}) \leq \underbrace{W_1(\overline{\mu_x}, \overline{\mu_{x,y}})}_A + \underbrace{W_1(\overline{\mu_{x,y}}, \overline{\mu_y})}_B. \quad (28)$$

*Term A.* We have  $W_1(\overline{\mu_x}, \overline{\mu_{x,y}}) \leq 2 \frac{|\mu_x| - |\mu_{x,y}|}{|\mu_x|} r$ . Hypothesis 5 still gives  $|\mu_x| \geq c_9 r^d$ . But Hypothesis 7 now yields

$$\begin{aligned} |\mu_x| - |\mu_{x,y}| &\leq \mu(\overline{\mathcal{B}}(x, r)) - \mu(\overline{\mathcal{B}}(x, r - \|x - y\|)) \\ &\leq c_{11} r^{d-\frac{1}{2}} \|x - y\|^{\frac{1}{2}}. \end{aligned}$$

We eventually obtain  $W_1(\overline{\mu_x}, \overline{\mu_{x,y}}) \leq 2 \frac{c_{11}}{c_9} r^{\frac{1}{2}} \|x - y\|^{\frac{1}{2}}$ .

*Term B.* In order to bound  $W_1(\overline{\mu_y}, \overline{\mu_{x,y}}) \leq 2 \frac{|\mu_y| - |\mu_{x,y}|}{|\mu_y|} r$ , Hypothesis 5 still gives  $|\mu_y| \geq (\frac{3}{4})^d c_9 r^d$ , and Hypothesis 7 yields

$$\begin{aligned} |\mu_y| - |\mu_{x,y}| &\leq \mu(\overline{\mathcal{B}}(y, r + \|x - y\|)) - \mu(\overline{\mathcal{B}}(y, r - \|x - y\|)) \\ &\leq c_{11} (r + \|x - y\|)^{d-\frac{1}{2}} (2\|x - y\|)^{\frac{1}{2}}, \end{aligned}$$

which is not greater than  $c_{11} (\frac{5}{4} r)^{d-\frac{1}{2}} (2\|x - y\|)^{\frac{1}{2}}$ . We finally get

$$W_1(\overline{\mu_y}, \overline{\mu_{x,y}}) \leq 2 \frac{c_{11} (\frac{5}{4} r)^{d-\frac{1}{2}} (2\|x - y\|)^{\frac{1}{2}}}{(\frac{3}{4})^d c_9 r^d} r \leq \frac{2^{\frac{5}{2}} 5^{d-\frac{1}{2}} c_{11}}{3^d c_9} r^{\frac{1}{2}} \|x - y\|^{\frac{1}{2}},$$

and we obtain the result by adding Terms A and B.  $\square$

We can now compare the localized measures of two probability measures.

**Lemma 4.7.** Let  $w = W_p(\mu, \nu)$ . Let  $y \in E$ . Suppose that there exists  $x \in \text{supp}(\mu)$  such that  $\|x - y\| \leq \alpha$  with  $\alpha = (\frac{w}{\rho^{d-1}})^{\frac{1}{2}}$ , and that  $\mu$  satisfies Hypotheses 5 and 6 at  $x$  with  $r < \min(\lambda(x), \frac{1}{2\rho})$ . Assume that  $w \leq \min(c_9, 1) (\frac{r}{4})^{d+1}$ . Then

$$W_1(\overline{\mu_y}, \overline{\nu_y}) \leq c_{17} \alpha,$$

$$\text{with } c_{17} = \frac{2^{d-1}}{c_9} + 2 \frac{12 \cdot 5^{d-1} c_{10} + 1}{3^d c_9} + 2^{d+3} \frac{(\frac{3}{2})^{d-1} c_{10} + 1}{c_9}.$$

*Proof.* Let  $\pi$  be an optimal transport for  $W_p(\mu, \nu)$ . Define  $\pi_y$  to be the restriction of the measure  $\pi$  to the set  $\overline{\mathcal{B}}(y, r) \times \overline{\mathcal{B}}(y, r) \subset E \times E$ . Its marginals  $p_{1*} \pi_y$  and  $p_{2*} \pi_y$  are submeasures of  $\mu_y$  and  $\nu_y$ . We shall use the triangle inequality:

$$W_1(\overline{\mu_y}, \overline{\nu_y}) \leq \underbrace{W_1(\overline{\mu_y}, \overline{p_{1*} \pi_y})}_A + \underbrace{W_1(\overline{p_{1*} \pi_y}, \overline{p_{2*} \pi_y})}_B + \underbrace{W_1(\overline{p_{2*} \pi_y}, \overline{\nu_y})}_C \quad (29)$$

Before examining each of these terms, note that we have

$$|\pi_y| = |p_{1*} \pi_y| = |p_{2*} \pi_y| \geq \mu(\overline{\mathcal{B}}(y, r - \alpha)) - \frac{w}{\alpha} \quad (30)$$



$$|v_y| \leq \mu(\overline{\mathcal{B}}(y, r + \alpha)) + \frac{w}{\alpha} \quad (31)$$

$$|v_y| \geq \mu(\overline{\mathcal{B}}(y, r - \alpha)) - \frac{w}{\alpha} \quad (32)$$

The first equation can be proven as follows:

$$\begin{aligned} \mu(\overline{\mathcal{B}}(y, r - \alpha)) &= \pi(\overline{\mathcal{B}}(y, r - \alpha) \times E) \\ &= \pi(\overline{\mathcal{B}}(y, r - \alpha) \times \overline{\mathcal{B}}(y, r)) + \pi(\overline{\mathcal{B}}(y, r - \alpha) \times \overline{\mathcal{B}}(y, r)^c) \end{aligned}$$

On the one hand,  $\pi(\overline{\mathcal{B}}(y, r - \alpha) \times \overline{\mathcal{B}}(y, r)) \leq \pi(\overline{\mathcal{B}}(y, r) \times \overline{\mathcal{B}}(y, r)) \leq |\pi_y|$ . On the other hand, Markov inequality yields

$$\pi(\overline{\mathcal{B}}(y, r - \alpha) \times \overline{\mathcal{B}}(y, r)^c) \leq \pi(\{(z, z'), \|z - z'\| \geq \alpha\}) \leq \frac{1}{\alpha} \int \|z - z'\| d\pi(z, z'),$$

and Jensen inequality gives

$$\frac{1}{\alpha} \int \|z - z'\| d\pi(z, z') \leq \frac{1}{\alpha} \left( \int \|z - z'\|^p d\pi(z, z') \right)^{\frac{1}{p}} = \frac{w}{\alpha}.$$

We deduce that  $\mu(\overline{\mathcal{B}}(y, r - \alpha)) \leq |\pi_y| + \frac{w}{\alpha}$ , which gives Equation (30). Equations (31) and (32) can be proven similarly. In addition to these preliminaries, note that the assumption  $w \leq \min(c_9, 1) \left(\frac{r}{4}\right)^{d+1}$  yields

$$\alpha \leq \frac{r}{4} \quad (33)$$

$$\frac{w}{\alpha} \leq \frac{c_9}{2} \left(\frac{r}{2}\right)^d \quad (34)$$

We now study the Terms B, A and C.

*Term B.* Since  $\overline{\pi_y} = \frac{\pi_y}{|\pi_y|}$  is a transport plan between  $\overline{p_{1*}\pi_y}$  and  $\overline{p_{2*}\pi_y}$ , we have

$$W_1(\overline{p_{1*}\pi_y}, \overline{p_{2*}\pi_y}) \leq \int \|z - z'\| \frac{d\pi_y(z, z')}{|\pi_y|} \leq \frac{1}{|\pi_y|} \int \|z - z'\| d\pi(z, z').$$

Moreover, Jensen inequality yields  $\int \|z - z'\| d\pi(z, z') \leq w$ . Hence

$$W_1(\overline{p_{1*}\pi_y}, \overline{p_{2*}\pi_y}) \leq \frac{w}{|\pi_y|}.$$

Let us prove that  $|\pi_y| \geq \frac{c_9}{2} \left(\frac{r}{2}\right)^d$ . According to Equation (30),  $|\pi_y| \geq \mu(\overline{\mathcal{B}}(y, r - \alpha)) - \frac{w}{\alpha}$ . Now, note that  $\mu(\overline{\mathcal{B}}(y, r - \alpha)) \geq \frac{c_9}{2^d} r^d$ . Indeed, using Hypothesis 5,

$$\mu(\overline{\mathcal{B}}(y, r - \alpha)) \geq \mu(\overline{\mathcal{B}}(x, r - \alpha - \|x - y\|)) \geq c_9(r - \alpha - \|x - y\|)^d,$$

and we conclude with  $\|x - y\| \leq \alpha \leq \frac{r}{4}$ . Now, using Equation (34), we get

$$\begin{aligned} |\pi_y| &\geq \mu(\overline{\mathcal{B}}(y, r - \alpha)) - \frac{w}{\alpha} \\ &\geq c_9 \left(\frac{r}{2}\right)^d - \frac{c_9}{2} \left(\frac{r}{2}\right)^d \geq \frac{c_9}{2} \left(\frac{r}{2}\right)^d. \end{aligned}$$

Finally, since  $\alpha = \left(\frac{w}{r^{d-1}}\right)^{\frac{1}{2}}$  and  $\alpha \leq \frac{r}{4}$ , we obtain

$$W_1(\overline{p_{1*}\pi_y}, \overline{p_{2*}\pi_y}) \leq \frac{w}{|\pi_y|} \leq \frac{w}{\frac{c_9}{2} \left(\frac{r}{2}\right)^d} = \frac{2^{d+1}}{c_9} \alpha^2 \frac{1}{r} \leq \frac{2^{d-1}}{c_9} \alpha.$$

*Term A.* According to Lemma 4.4, we have

$$W_1(\overline{\mu_y}, \overline{p_{1*}\pi_y}) \leq 2 \frac{|\mu_y| - |p_{1*}\pi_y|}{|\mu_y|} r. \quad (35)$$

We can use Equation (30) to get

$$\begin{aligned} |\mu_y| - |p_{1*}\pi_y| &\leq \mu(\overline{\mathcal{B}}(y, r)) - \mu(\overline{\mathcal{B}}(y, r - \alpha)) + \frac{w}{\alpha} \\ &\leq \mu(\overline{\mathcal{B}}(x, r + \|x - y\|)) - \mu(\overline{\mathcal{B}}(x, r - \alpha - \|x - y\|)) + \frac{w}{\alpha}. \end{aligned}$$

Moreover, by Hypothesis 6, we have

$$\mu(\overline{\mathcal{B}}(x, r + \|x - y\|)) - \mu(\overline{\mathcal{B}}(x, r - \alpha - \|x - y\|)) \leq c_{10}(r + \|x - y\|)^{d-1}(2\|x - y\| + \alpha),$$

which is not greater than  $c_{10}(\frac{5}{4}r)^{d-1}3\alpha$  since  $\|x - y\| \leq \alpha \leq \frac{r}{4}$ . Besides,  $\frac{w}{\alpha} = r^{d-1}\alpha$ , and we obtain

$$|\mu_y| - |p_{1*}\pi_y| \leq \left(3 \left(\frac{5}{4}\right)^{d-1} c_{10} + 1\right) r^{d-1} \alpha.$$

Finally, thanks to Hypothesis 5, we write

$$\begin{aligned} |\mu_y| = \mu(\overline{\mathcal{B}}(y, r)) &\geq \mu(\overline{\mathcal{B}}(x, r - \|x - y\|)) \\ &\geq c_9(r - \|x - y\|)^d \geq c_9 \left(\frac{3}{4}\right)^d r^d \end{aligned}$$

and we obtain

$$\frac{|\mu_y| - |p_{1*}\pi_y|}{|\mu_y|} \leq \frac{((3(\frac{5}{4})^{d-1} c_{10} + 1) r^{d-1})}{c_9(\frac{3}{4})^d r^d} \alpha = \frac{1}{r} \cdot \frac{12 \cdot 5^{d-1} c_{10} + 1}{3^d c_9} \alpha.$$

Combined with Equation (35), we deduce

$$W_1(\overline{\mu_y}, \overline{p_{1*}\pi_y}) \leq 2 \frac{12 \cdot 5^{d-1} c_{10} + 1}{3^d c_9} \alpha.$$

*Term C.* It is similar to Term A. First, one shows that

$$W_1(\overline{v_y}, \overline{p_{2*}\pi_y}) \leq 2 \frac{|v_y| - |p_{2*}\pi_y|}{|v_y|} r. \quad (36)$$

Using Equations (30) and (31) we get

$$\begin{aligned} |v_y| - |p_{2*}\pi_y| &\leq \mu(\overline{\mathcal{B}}(y, r + \alpha)) + \frac{w}{\alpha} - \mu(\overline{\mathcal{B}}(y, r - \alpha)) + \frac{w}{\alpha} \\ &\leq \mu(\overline{\mathcal{B}}(x, r + \|x - y\| + \alpha)) - \mu(\overline{\mathcal{B}}(x, r - \alpha - \|x - y\|)) + 2\frac{w}{\alpha}. \end{aligned}$$

By Hypothesis 6, we have

$$\begin{aligned} &\mu(\overline{\mathcal{B}}(x, r + \|x - y\| + \alpha)) - \mu(\overline{\mathcal{B}}(x, r - \alpha - \|x - y\|)) \\ &\leq c_{10}(r + \|x - y\| + \alpha)^{d-1}(2\|x - y\| + 2\alpha) \end{aligned}$$

which is not greater than  $c_{10}(\frac{3}{2}r)^{d-1}4\alpha$  since  $\|x - y\| \leq \alpha \leq \frac{r}{4}$ . Moreover,  $\frac{w}{\alpha} = r^{d-1}\alpha$ , and we obtain

$$|v_y| - |p_{2*}\pi_y| \leq (4(\frac{3}{2})^{d-1} c_{10} + 2) r^{d-1} \alpha.$$

We have seen that

$$|v_y| \geq \mu(\overline{\mathcal{B}}(y, r - \alpha)) - \frac{w}{\alpha} \geq \frac{c_9}{2} \left(\frac{r}{2}\right)^d.$$

Hence

$$\frac{|v_y| - |p_{2*}\pi_y|}{|v_y|} \leq \frac{(4(\frac{3}{2})^{d-1}c_{10} + 2)r^{d-1}}{\frac{c_9}{2}(\frac{r}{2})^d} \alpha = \frac{1}{r} \cdot 2^{d+2} \frac{(\frac{3}{2})^{d-1}b + 1}{c_9} \alpha,$$

and we deduce from Equation (36) that

$$W_1(\overline{\mu_y}, \overline{p_{1*}\pi_y}) \leq 2^{d+3} \frac{(\frac{3}{2})^{d-1}c_{10} + 1}{c_9} \alpha.$$

To conclude, summing up the Terms A, B and C gives  $W_1(\overline{\mu_y}, \overline{v_y}) \leq c_{17}\alpha$  with

$$c_{17} = \frac{2^{d-1}}{c_9} + 2 \frac{12 \cdot 5^{d-1}c_{10} + 1}{3^d c_9} + 2^{d+3} \frac{(\frac{3}{2})^{d-1}c_{10} + 1}{c_9},$$

as wanted.  $\square$

As before, we prove a version of Lemma 4.7 where Hypothesis 6 is replaced by the weaker Hypothesis 7.

**Lemma 4.8.** *Let  $w = W_p(\mu, \nu)$ . Let  $y \in E$ . Suppose that there exists  $x \in \text{supp}(\mu)$  such that  $\|x - y\| \leq \alpha$  with  $\alpha = (\frac{w}{r^{d-1}})^{\frac{1}{2}}$ , and that  $\mu$  satisfies Hypotheses 5 and 7 at  $x$  with  $r < \frac{1}{2\rho}$ . Assume that  $w \leq \min(c_9, 1)(\frac{r}{4})^{d+1}$ . Then*

$$W_1(\overline{\mu_y}, \overline{v_y}) \leq c_{18} r^{\frac{1}{2}} \alpha^{\frac{1}{2}},$$

$$\text{with } c_{18} = \frac{2^{d-2}}{c_9} + \frac{4 \cdot 3^{\frac{1}{2}} 5^{d-\frac{1}{2}} c_{11} + 4^{d-\frac{1}{2}}}{3^d c_9} + 2 \cdot 4^d \frac{2c_{11}(\frac{3}{2})^{d-\frac{1}{2}} + 1}{3^d c_9}.$$

*Proof.* The proof is similar as Lemma 4.7. Let us highlight the modifications. Since  $\alpha \leq \frac{r}{4}$  and  $\frac{w}{\alpha} = r^{d-1}\alpha$ , we have the inequalities

$$\alpha^{\frac{1}{2}} \leq \frac{1}{2} r^{\frac{1}{2}} \quad \text{and} \quad \frac{w}{\alpha} \leq \frac{1}{2} r^{d-\frac{1}{2}} \alpha^{\frac{1}{2}}.$$

We still write the triangle inequality:

$$W_1(\overline{\mu_y}, \overline{v_y}) \leq \underbrace{W_1(\overline{\mu_y}, \overline{p_{1*}\pi_y})}_A + \underbrace{W_1(\overline{p_{1*}\pi_y}, \overline{p_{2*}\pi_y})}_B + \underbrace{W_1(\overline{p_{2*}\pi_y}, \overline{v_y})}_C \quad (37)$$

where  $\pi$  is an optimal transport plan for  $W_p(\mu, \nu)$ .

*Term B.* The argument to obtain  $W_1(\overline{p_{1*}\pi_y}, \overline{p_{2*}\pi_y}) \leq \frac{2^{d-1}}{c_9} \alpha$  is unchanged, and we use  $\alpha^{\frac{1}{2}} \leq \frac{1}{2} r^{\frac{1}{2}}$  to get

$$W_1(\overline{p_{1*}\pi_y}, \overline{p_{2*}\pi_y}) \leq \frac{2^{d-2}}{c_9} \alpha^{\frac{1}{2}} r^{\frac{1}{2}}.$$

*Term A.* Using Hypothesis 7, we have

$$\begin{aligned} & \mu(\overline{\mathcal{B}}(x, r + \|x - y\|)) - \mu(\overline{\mathcal{B}}(x, r - \alpha - \|x - y\|)) \\ & \leq c_{11}(r + \|x - y\|)^{d-\frac{1}{2}} (2\|x - y\| + \alpha)^{\frac{1}{2}} \\ & \leq c_{11} \left(\frac{5}{4}r\right)^{d-\frac{1}{2}} 3^{\frac{1}{2}} \alpha^{\frac{1}{2}}. \end{aligned}$$

And since  $\frac{w}{\alpha} \leq \frac{1}{2} r^{d-\frac{1}{2}} \alpha^{\frac{1}{2}}$ , we get

$$\begin{aligned} |\mu_y| - |p_{1*}\pi_y| &\leq \mu(\overline{\mathcal{B}}(x, r + \|x-y\|)) - \mu(\overline{\mathcal{B}}(x, r - \alpha - \|x-y\|)) + \frac{w}{\alpha} \\ &\leq \left( c_{11} \left( \frac{5}{4} \right)^{d-\frac{1}{2}} 3^{\frac{1}{2}} + \frac{1}{2} \right) r^{d-\frac{1}{2}} \alpha^{\frac{1}{2}}. \end{aligned}$$

Finally, we use

$$\begin{aligned} |\mu_y| &= \mu(\overline{\mathcal{B}}(y, r)) \geq \mu(\overline{\mathcal{B}}(x, r - \|x-y\|)) \\ &\geq c_9 (r - \|x-y\|)^d \geq c_9 \left( \frac{3}{4} \right)^d r^d \end{aligned}$$

to obtain

$$\frac{|\mu_y| - |p_{1*}\pi_y|}{|\mu_y|} \leq \frac{((c_{11}(\frac{5}{4})^{d-\frac{1}{2}} 3^{\frac{1}{2}} + \frac{1}{2}) r^{d-\frac{1}{2}}) \alpha^{\frac{1}{2}}}{c_9 (\frac{3}{4})^d r^d} \alpha^{\frac{1}{2}} = \frac{1}{r^{\frac{1}{2}}} \cdot \frac{2 \cdot 3^{\frac{1}{2}} 5^{d-\frac{1}{2}} c_{11} + 4^{d-\frac{1}{2}}}{3^d c_9} \alpha^{\frac{1}{2}}$$

and we deduce

$$W_1(\overline{\mu_y}, \overline{p_{1*}\pi_y}) \leq 2 \frac{|\mu_y| - |p_{1*}\pi_y|}{|\mu_y|} r \leq \frac{4 \cdot 3^{\frac{1}{2}} 5^{d-\frac{1}{2}} c_{11} + 4^{d-\frac{1}{2}}}{3^d c_9} r^{\frac{1}{2}} \alpha^{\frac{1}{2}}.$$

*Term C.* We use Hypothesis 7 to get

$$\begin{aligned} &\mu(\overline{\mathcal{B}}(x, r + \|x-y\| + \alpha)) - \mu(\overline{\mathcal{B}}(x, r - \alpha - \|x-y\|)) \\ &\leq c_{11} (r + \|x-y\| + \alpha)^{d-\frac{1}{2}} (2\|x-y\| + 2\alpha)^{\frac{1}{2}} \\ &\leq 2c_{11} \left( \frac{3}{2} r \right)^{d-\frac{1}{2}} \alpha^{\frac{1}{2}}. \end{aligned}$$

And since  $\frac{w}{\alpha} \leq \frac{1}{2} r^{d-\frac{1}{2}} \alpha^{\frac{1}{2}}$ , we get

$$\begin{aligned} |v_y| - |p_{2*}\pi_y| &\leq \mu(\overline{\mathcal{B}}(x, r + \|x-y\| + \alpha)) - \mu(\overline{\mathcal{B}}(x, r - \alpha - \|x-y\|)) + 2\frac{w}{\alpha} \\ &\leq \left( 2c_{11} \left( \frac{3}{2} \right)^{d-\frac{1}{2}} + 1 \right) r^{d-\frac{1}{2}} \alpha^{\frac{1}{2}}. \end{aligned}$$

Finally, we use

$$\begin{aligned} |\mu_y| &= \mu(\overline{\mathcal{B}}(y, r)) \geq \mu(\overline{\mathcal{B}}(x, r - \|x-y\|)) \\ &\geq c_9 (r - \|x-y\|)^d \geq c_9 \left( \frac{3}{4} \right)^d r^d \end{aligned}$$

to obtain

$$\frac{|\mu_y| - |p_{1*}\pi_y|}{|\mu_y|} \leq \frac{(2c_{11}(\frac{3}{2})^{d-\frac{1}{2}} + 1) r^{d-\frac{1}{2}} \alpha^{\frac{1}{2}}}{c_9 (\frac{3}{4})^d r^d} \alpha^{\frac{1}{2}} = \frac{1}{r^{\frac{1}{2}}} \cdot 4^d \frac{2c_{11}(\frac{3}{2})^{d-\frac{1}{2}} + 1}{3^d c_9} \alpha^{\frac{1}{2}}$$

and we deduce

$$W_1(\overline{\mu_y}, \overline{p_{1*}\pi_y}) \leq 2 \frac{|\mu_y| - |p_{1*}\pi_y|}{|\mu_y|} r \leq 2 \cdot 4^d \frac{2c_{11}(\frac{3}{2})^{d-\frac{1}{2}} + 1}{3^d c_9} r^{\frac{1}{2}} \alpha^{\frac{1}{2}}.$$

We finally obtain the result by summing Terms A, B and C.  $\square$

*Remark 4.9.* Let us comment the inequality of Lemma 4.7 with  $p = 1$ , valid for all  $r$  such that  $W_1(\mu, \nu) \leq \min(a, 1) \left(\frac{r}{4}\right)^{d+1}$ :

$$W_1(\overline{\mu}_y, \overline{\nu}_y) \leq c_{17} \left( \frac{W_1(\mu, \nu)}{r^{d-1}} \right)^{\frac{1}{2}}. \quad (38)$$

If  $r$  is assumed to be constant, the behavior of  $W_1(\overline{\mu}_y, \overline{\nu}_y)$ , when  $W_1(\mu, \nu)$  goes to 0, is  $W_1(\overline{\mu}_y, \overline{\nu}_y) \lesssim W_1(\mu, \nu)^{\frac{1}{2}}$ . On the other hand, if  $r$  is supposed to follow the worst case, i.e.,  $r$  is of order  $W_1(\mu, \nu)^{\frac{1}{d+1}}$ , then  $W_1(\overline{\mu}_y, \overline{\nu}_y)$  is of order  $W_1(\overline{\mu}_y, \overline{\nu}_y) \lesssim W_1(\mu, \nu)^{\frac{1}{d+1}}$ .

A similar stability result already appears in [39, Theorem 4.3], where  $\mu$  and  $\nu$  are two probability measures on a bounded set  $X$ , and  $\mu$  satisfy the following condition:  $\forall x \in X, \forall s, r \leq 0$  s.t.  $s \leq r$ , we have  $\frac{\mu(\overline{\mathcal{B}}(x, r))}{\mu(\overline{\mathcal{B}}(x, s))} \leq \left(\frac{r}{s}\right)^d$ . The theorem states that, denoting  $D = \text{diam}(X)$ , for all  $x \in X$ ,

$$W_1(\overline{\mu}_x, \overline{\nu}_x) \leq (1 + 2r) \left[ \frac{W_1(\mu, \nu)^{\frac{1}{2}}}{\min(1, (\frac{r}{D})^d)} + \left(1 + \frac{W_1(\mu, \nu)^{\frac{1}{2}}}{r}\right)^d - 1 \right].$$

When  $r \leq D$  and  $W_1(\mu, \nu)$  goes to zero, we obtain that  $W_1(\overline{\mu}_x, \overline{\nu}_x)$  is of order

$$W_1(\overline{\mu}_x, \overline{\nu}_x) \leq (1 + 2r) D^d \left( \frac{W_1(\mu, \nu)}{r^{2d}} \right)^{\frac{1}{2}}.$$

The exponent on  $r$  is greater here than in Equation (38).

Let us show that the order we obtained,  $\left(\frac{W_1(\mu, \nu)}{r^{d-1}}\right)^{\frac{1}{2}}$ , is optimal. More precisely, let us show that, for every  $d \geq 1$ ,  $r > 0$  and  $\varepsilon > 0$  fixed, there exists measures  $\mu$  and  $\nu$  on  $\mathbb{R}^d$  that satisfies the assumptions of Lemma 4.7, but such that

$$W_1(\overline{\mu}_y, \overline{\nu}_y) \geq c_d \left( \frac{W_1(\mu, \nu)}{r^{d-1}} \right)^{\frac{1}{2}} - \varepsilon$$

with  $c_d = \frac{1}{d+1} \left( \frac{2d}{V_d} \right)^{\frac{1}{2}}$ . We consider the following example. Let  $\mu = \mathcal{H}_{[0,1]^d}^d$  be the Lebesgue measure on the hypercube  $[0, 1]^d$ . Denote  $y = (\frac{1}{2}, \dots, \frac{1}{2})$  its center,  $B = \mathcal{B}(y, r)$  the open ball, and  $A$  the annulus defined as

$$A = \mathcal{B}(y, r + \varepsilon) \setminus \mathcal{B}(y, r)$$

where  $0 < \varepsilon < r < \frac{1}{4}$ . In the following,  $r$  stays fixed, and  $\varepsilon$  will go to zero. Consider the probability measure

$$\nu = \mathcal{H}_{[0,1]^d \setminus A}^d + \frac{V_d(r + \varepsilon)^d - V_d r^d}{S_{d-1} r^{d-1}} \mathcal{H}_{\partial \mathcal{B}(y, r)}^{d-1}.$$

Let  $\overline{\mu}_y$  and  $\overline{\nu}_y$  be the localized probability measures associated to  $\mu$  and  $\nu$  with parameter  $r$ . We shall show that

$$W_1(\mu, \nu) \text{ is of order } r^{d-1} \varepsilon^2 \quad \text{and} \quad W_1(\overline{\mu}_y, \overline{\nu}_y) \text{ is of order } \varepsilon$$

when  $\varepsilon \rightarrow 0$ . These measures are depicted in Figure 21.

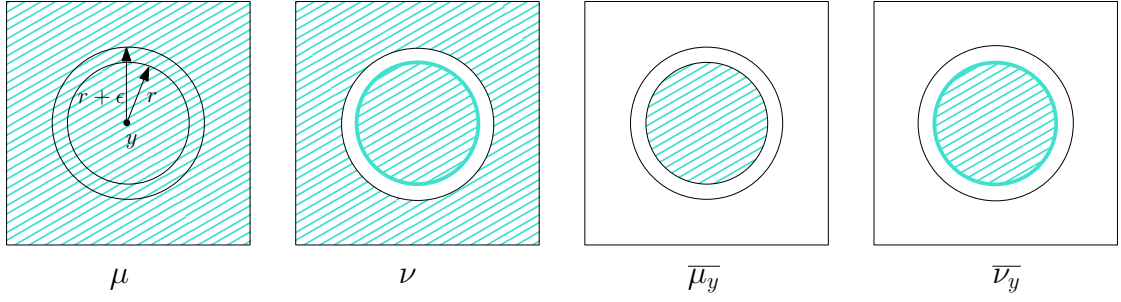


Figure 21: The measures involved in the example. A hatched area represents the  $d$ -dimensional Hausdorff measure  $\mathcal{H}^d$ , and a bold circle represents the  $(d-1)$ -dimensional Hausdorff measure  $\mathcal{H}^{d-1}$ .

*Step 1: Study of  $W_1(\mu, \nu)$ .* An optimal transport plan between  $\mu$  and  $\nu$  is given by transporting the submeasure  $\mathcal{H}_A^d$  of  $\mu$  onto the submeasure  $\frac{V_d(r+\varepsilon)^d - V_d r^d}{S_{d-1} r^{d-1}} \mathcal{H}_{\partial \mathcal{B}(y, r)}^{d-1}$  of  $\nu$  via the map  $x \mapsto \frac{r}{\|x\|} x$ . Consequently, the Wasserstein distance is

$$W_1(\mu, \nu) = \int_A \left\| x - \frac{r}{\|x\|} x \right\| \frac{V_d(r+\varepsilon)^d - V_d r^d}{S_{d-1} r^{d-1}} d\mathcal{H}^d(x).$$

A change of coordinates shows that

$$\int_A \left\| x - \frac{r}{\|x\|} x \right\| d\mathcal{H}^d(x) = \int_{\partial \mathcal{B}(0,1)} \int_r^{r+\varepsilon} (t-r) t^{d-1} d\mathcal{H}^1(t) d\mathcal{H}^{d-1}(v).$$

Let us split the integral as

$$\int_r^{r+\varepsilon} (t-r) t^{d-1} d\mathcal{H}^1(t) = \int_r^{r+\varepsilon} t^d d\mathcal{H}^1(t) - \int_r^{r+\varepsilon} r t^{d-1} d\mathcal{H}^1(t).$$

On the one hand, we have

$$\begin{aligned} \int_r^{r+\varepsilon} t^d d\mathcal{H}^1(t) &= \frac{1}{d+1} \left( (r+\varepsilon)^{d+1} - r^{d+1} \right) \\ &= r^d \varepsilon + \frac{d}{2} r^{d-1} \varepsilon^2 + o(\varepsilon^2), \end{aligned}$$

where the Little-O notation refers to  $\varepsilon \rightarrow 0$ . On the other hand,

$$\begin{aligned} \int_r^{r+\varepsilon} r t^{d-1} d\mathcal{H}^1(t) &= r \left( r^{d-1} \varepsilon + \frac{d-1}{2} r^{d-2} \varepsilon^2 + o(\varepsilon^2) \right) \\ &= r^d \varepsilon + \frac{d-1}{2} r^{d-1} \varepsilon^2 + o(\varepsilon^2). \end{aligned}$$

We deduce that  $\int_r^{r+\varepsilon} (t-r) t^{d-1} d\mathcal{H}^1(t) = \frac{1}{2} r^{d-1} \varepsilon^2 + o(\varepsilon^2)$ , and

$$\int_A \left\| x - \frac{r}{\|x\|} x \right\| d\mathcal{H}^d(x) = \frac{S_{d-1}}{2} r^{d-1} \varepsilon^2 + o(\varepsilon^2).$$

In other words,  $W_1(\mu, \nu) = \frac{dV_d}{2} r^{d-1} \varepsilon^2 + o(\varepsilon^2)$ .

*Step 2: Study of  $W_1(\overline{\mu}_y, \overline{\nu}_y)$ .* Consider the measures

$$\begin{aligned} \overline{\mu}_x &= \frac{1}{V_d r^d} \mathcal{H}_B^d = \left( \frac{1}{V_d (r+\varepsilon)^d} + \frac{V_d (r+\varepsilon)^d - V_d r^d}{V_d (r+\varepsilon)^d V_d r^d} \right) \mathcal{H}_B^d, \\ \overline{\nu}_x &= \frac{1}{V_d (r+\varepsilon)^d} \left( \mathcal{H}_B^d + \frac{V_d (r+\varepsilon)^d - V_d r^d}{S_{d-1} r^{d-1}} \mathcal{H}_{\partial \mathcal{B}(y, r)}^{d-1} \right). \end{aligned}$$

Let us compute the Wasserstein distance  $W_1(\overline{\mu}_y, \overline{\nu}_y)$ . As before, an optimal transport plan is given by transporting the submeasure  $\frac{V_d(r+\varepsilon)^d - V_d r^d}{V_d(r+\varepsilon)^d V_d r^d} \mathcal{H}_B^d$  of  $\overline{\mu}_x$  onto the submeasure  $\frac{V_d(r+\varepsilon)^d - V_d r^d}{V_d(r+\varepsilon)^d S_{d-1} r^{d-1}} \mathcal{H}_{\partial \mathcal{B}(y,r)}^{d-1}$  of  $\overline{\nu}_x$ . We have:

$$W_1(\overline{\mu}_y, \overline{\nu}_y) = \int_B \left\| x - \frac{r}{\|x\|} x \right\| \frac{V_d(r+\varepsilon)^d - V_d r^d}{V_d(r+\varepsilon)^d V_d r^d} d\mathcal{H}^d(x)$$

A change of coordinates yields

$$\int_B \left\| x - \frac{r}{\|x\|} x \right\| d\mathcal{H}^d(x) = \frac{S_{d-1}}{d(d+1)} r^{d+1}.$$

Besides, we have

$$\frac{V_d(r+\varepsilon)^d - V_d r^d}{V_d(r+\varepsilon)^d V_d r^d} = \frac{dV_d r^{d-1} \varepsilon + O(\varepsilon^2)}{V_d(r+\varepsilon)^d V_d r^d} = \frac{d}{V_d} \frac{\varepsilon}{r^{d+1}} + O(\varepsilon^2).$$

We deduce that

$$W_1(\overline{\mu}_y, \overline{\nu}_y) = \frac{S_{d-1}}{d(d+1)} \frac{d}{V_d} \varepsilon + O(\varepsilon^2) = \frac{d}{d+1} \varepsilon + O(\varepsilon^2).$$

*Step 3: Comparison of the distances.* Using  $W_1(\mu, \nu) = \frac{dV_d}{2} r^{d-1} \varepsilon^2 + o(\varepsilon^2)$  and  $W_1(\overline{\mu}_y, \overline{\nu}_y) = \frac{d}{d+1} \varepsilon + O(\varepsilon^2)$ , we get

$$\frac{W_1(\overline{\mu}_y, \overline{\nu}_y)^2}{W_1(\mu, \nu)} = c_d \frac{1}{r^{d-1}} + O(\varepsilon) \quad \text{where} \quad c_d = \frac{\left(\frac{d}{d+1}\right)^2}{\frac{dV_d}{2}} = \frac{2d}{(d+1)^2 V_d}.$$

Together with  $W_1(\mu, \nu)^{\frac{1}{2}} = O(\varepsilon)$ , we deduce the result:

$$W_1(\overline{\mu}_y, \overline{\nu}_y) = c^{\frac{1}{2}} \left( \frac{W_1(\mu, \nu)}{r^{d-1}} \right)^{\frac{1}{2}} + O(\varepsilon^2).$$

#### 4.4 Stability of the estimation

In this subsection we study the stability of the normalized local covariance matrix operator  $\mu \mapsto \overline{\Sigma}_\mu(\cdot)$  (see Definition 4.1) with respect to the  $W_p$  metric on measures.

As an introduction to the problem, let  $\mu$  and  $\nu$  be two probability measures,  $x \in \text{supp}(\mu) \cap \text{supp}(\nu)$ , and consider the Frobenius distance  $\|\overline{\Sigma}_\mu(x) - \overline{\Sigma}_\nu(x)\|_F$  between the normalized local covariance matrices. One shows that this distance is related to the 1-Wasserstein distance between the localized probability measures  $\overline{\mu}_x$  and  $\overline{\nu}_x$  via the following inequality (see Equation (41) in the proof of Lemma 4.12):

$$\|\overline{\Sigma}_\mu(x) - \overline{\Sigma}_\nu(x)\|_F \leq \frac{2}{r} W_1(\overline{\mu}_x, \overline{\nu}_x).$$

Without any assumption on the measures, it is not true that  $W_1(\overline{\mu}_x, \overline{\nu}_x)$  goes to 0 as  $W_1(\mu, \nu)$  does (see Remark 4.10). However, if we assume that  $\mu$  satisfies the Hypotheses 5 and 6, that  $x$  satisfies  $\lambda(x) > 0$  and that  $r$  is chosen such that

$$4 \left( \frac{W_1(\mu, \nu)}{\min(c_9, 1)} \right)^{\frac{1}{d+1}} \leq r < \min \left( \lambda(x), \frac{1}{2\rho} \right),$$

then we have seen in Lemma 4.7 that

$$W_1(\overline{\mu}_x, \overline{\nu}_x) \leq c_{17} \left( \frac{W_1(\mu, \nu)}{r^{d-1}} \right)^{\frac{1}{2}}. \quad (39)$$

As a consequence of this inequality, estimating local covariance matrices is robust in Wasserstein distance:

$$\|\bar{\Sigma}_\mu(x) - \bar{\Sigma}_\nu(x)\|_F \leq 2c_{17} \left( \frac{W_1(\mu, \nu)}{r^{d+1}} \right)^{\frac{1}{2}}. \quad (40)$$

We point out that another result of this kind bounds the distance  $\|\bar{\Sigma}_\mu(x) - \bar{\Sigma}_\nu(x)\|_F$  with the  $\infty$ -Wasserstein distance  $W_\infty(\mu, \nu)$  [38, Theorem 3]. Namely, if  $\mu$  and  $\nu$  are fully supported probability measures with densities upper bounded by  $l > 0$  and supports included in  $X \subset \mathbb{R}^d$ , denoting  $D = \text{diam}(X)$ , we have

$$\|\bar{\Sigma}_\mu(x) - \bar{\Sigma}_\nu(x)\|_F \leq l A W_\infty(\mu, \nu),$$

$$\text{where } A = \frac{d}{d+2} \frac{(r+D)^{d+1}}{D r^d} + \frac{(2r+D)(r+D)^d}{r^d} + \frac{2d}{d+2} \frac{(r+D)^{d+2}}{D r^d}.$$

*Remark 4.10.* In general, for  $x \in \text{supp}(\mu) \cap \text{supp}(\nu)$ , it is not true that  $\|\bar{\Sigma}_\mu(x) - \bar{\Sigma}_\nu(x)\|_F$  goes to zero as  $W_1(\mu, \nu)$  goes to zero. To see this, one can consider  $\varepsilon > 0$ , and the measures on  $\mathbb{R}$  defined as

$$\mu = \frac{1}{2}(\delta_0 + \delta_1) \quad \text{and} \quad \nu = \frac{1}{2}(\delta_0 + \delta_{1+\varepsilon}),$$

where  $\delta_x$  denotes the Dirac mass on  $x$ . Choose the scale parameter  $r = 1$ . Restricting the measures  $\mu$  and  $\nu$  to the ball  $\mathcal{B}(0, 1)$  of  $\mathbb{R}$  gives  $\bar{\mu}_0 = \frac{1}{2}(\delta_0 + \delta_1)$  and  $\bar{\nu}_0 = \delta_0$ . According to Definition 4.1, we deduce that the local covariance matrices are

$$\Sigma_\mu(0) = \frac{1}{2}1^{\otimes 2} \quad \text{and} \quad \Sigma_\nu(0) = 0.$$

Hence  $\|\bar{\Sigma}_\mu(x) - \bar{\Sigma}_\nu(x)\|_F = \frac{1}{2}$ , independently of  $\varepsilon$ . But, on the other hand, we have  $W_1(\mu, \nu) = \frac{1}{2}\varepsilon$ , which goes to zero.

Similarly,  $W_{p,\gamma}(\check{\mu}, \check{\nu})$  does not have to go to zero when  $W_1(\mu, \nu)$  does. Indeed, a similar computation shows that the local covariance matrices at 1 are

$$\Sigma_\mu(1) = \frac{1}{2}1^{\otimes 2} \quad \text{and} \quad \Sigma_\nu(1) = 0,$$

and we deduce that the lifted measures are

$$\check{\mu} = \frac{1}{2} \left( \delta_{(0, \frac{1}{2}1^{\otimes 2})} + \delta_{(1, \frac{1}{2}1^{\otimes 2})} \right) \quad \text{and} \quad \check{\nu} = \frac{1}{2} \left( \delta_{(0,0)} + \delta_{(1+\varepsilon,0)} \right).$$

Using the optimal transport plan  $\pi$  between  $\check{\mu}$  and  $\check{\nu}$  that sends  $\delta_{(0, \frac{1}{2}1^{\otimes 2})}$  to  $\delta_{(0,0)}$  and  $\delta_{(1, \frac{1}{2}1^{\otimes 2})}$  to  $\delta_{(1+\varepsilon,0)}$ , we get

$$\begin{aligned} W_{p,\gamma}^p(\check{\mu}, \check{\nu}) &= \frac{1}{2} \left\| \left( 0, \frac{1}{2}1^{\otimes 2} \right) - (0,0) \right\|_\gamma^p + \frac{1}{2} \left\| \left( 1, \frac{1}{2}1^{\otimes 2} \right) - (1+\varepsilon,0) \right\|_\gamma^p \\ &= \frac{1}{2} \left( \left( \frac{\gamma}{2} \right)^p + \left( \varepsilon^2 + \gamma^2 \frac{1}{4} \right)^{\frac{p}{2}} \right) \geq \left( \frac{\gamma}{2} \right)^p. \end{aligned}$$

Again, we see that  $W_{p,\gamma}^p(\check{\mu}, \check{\nu})$  is lower bounded, independently of  $\varepsilon$ . Hence  $W_{p,\gamma}(\check{\mu}, \check{\nu})$  does not go to zero as  $W_1(\mu, \nu)$  does. However, under regularity assumptions on  $\mu$ , the following proposition states that it is the case.

**Proposition 4.11.** *Let  $\mu$  and  $\nu$  be two probability measures on  $E$ . Suppose that  $\mu$  satisfies Hypotheses 5, 6 and 7. Define  $w = W_p(\mu, \nu)$ . Suppose that  $r \leq \min\left(\frac{1}{2p}, 1\right)$  and  $w \leq \min(c_9, 1) \left(\frac{r}{4}\right)^{d+1}$ . Then*

$$W_{p,\gamma}(\check{\mu}, \check{\nu}) \leq 2w + \gamma c_{19} \left( \frac{w}{r^{d+1}} \right)^{\frac{1}{2}} + \gamma c'_{19} \mu(\lambda^r)^{\frac{1}{p}} \left( \frac{w}{r^{d+1}} \right)^{\frac{1}{4}},$$

with  $c_{19} = 4(1 + c_{20})$  and  $c'_{19} = 4c_{18}$ .



*Proof.* According to Lemma 4.12 stated below, we have

$$W_{p,\gamma}(\check{\mu}, \check{\nu}) \leq 2^{\frac{p-1}{p}} \left(1 + \frac{2\gamma}{r}\right) w + 2^{\frac{p-1}{p}} \frac{2\gamma}{r} \left( \int W_1^p(\overline{\mu}_x, \overline{\nu}_y) d\pi(x, y) \right)^{\frac{1}{p}}.$$

Let  $\alpha = \left(\frac{w}{r^{d-1}}\right)^{\frac{1}{2}}$ . Lemma 4.13, also stated below, gives

$$\left( \int W_1(\overline{\mu}_x, \overline{\nu}_y) d\pi(x, y) \right)^{\frac{1}{p}} \leq 2^{\frac{p-1}{p}} \left( c_{18} r^{\frac{1}{2}} \mu(\lambda^r)^{\frac{1}{p}} \alpha^{\frac{1}{2}} + c_{20} \alpha \right)$$

Combining these inequalities yields

$$\begin{aligned} W_{p,\gamma}(\check{\mu}, \check{\nu}) &\leq 2^{\frac{p-1}{p}} w + 2^{\frac{p-1}{p}} \frac{2\gamma}{r} \left( w + 2^{\frac{p-1}{p}} c_{20} \alpha \right) + \left( 2^{\frac{p-1}{p}} \right)^2 \frac{2\gamma}{r} c_{18} r^{\frac{1}{2}} \mu(\lambda^r)^{\frac{1}{p}} \alpha^{\frac{1}{2}} \\ &\leq 2w + 2 \cdot 2\gamma \left( \frac{w}{r} + 2c_{20} \frac{\alpha}{r} \right) + 2^2 \cdot 2\gamma c_{18} \mu(\lambda^r)^{\frac{1}{p}} \left( \frac{\alpha}{r} \right)^{\frac{1}{2}}, \end{aligned}$$

where we used  $2^{\frac{p-1}{p}} \leq 2$ . Beside, since  $r \leq 1$ , we have  $w \leq 1$ , and

$$w = \left( \frac{w}{r^{d-1}} \right)^{\frac{1}{2}} r^{\frac{d-1}{2}} w^{\frac{1}{2}} \leq \left( \frac{w}{r^{d-1}} \right)^{\frac{1}{2}} = \alpha.$$

Consequently, we have

$$W_{p,\gamma}(\check{\mu}, \check{\nu}) \leq 2^{\frac{p-1}{p}} w + 2^{\frac{p-1}{p}} 2\gamma \left( 1 + 2^{\frac{p-1}{p}} c_{20} \right) \frac{\alpha}{r} + \left( 2^{\frac{p-1}{p}} \right)^2 2\gamma c_{18} \mu(\lambda^r)^{\frac{1}{p}} \left( \frac{\alpha}{r} \right)^{\frac{1}{2}}.$$

We obtain the result by replacing  $\frac{\alpha}{r}$  with  $\left(\frac{w}{r^{d+1}}\right)^{\frac{1}{2}}$ . □

Let us interpret the inequality given by Proposition 4.11:

$$W_{p,\gamma}(\check{\mu}, \check{\nu}) \leq 2W_p(\mu, \nu) + \gamma c_{19} \left( \frac{W_p(\mu, \nu)}{r^{d+1}} \right)^{\frac{1}{2}} + \gamma c'_{19} \mu(\lambda^r)^{\frac{1}{p}} \left( \frac{W_p(\mu, \nu)}{r^{d+1}} \right)^{\frac{1}{4}},$$

The lifted measures  $\check{\mu}$  and  $\check{\nu}$  are defined on the lift space  $E \times M(E)$ . Hence the Wasserstein distance  $W_{p,\gamma}(\check{\mu}, \check{\nu})$  may witness a difference with respect to the Euclidean coordinate ( $E$ -coordinate) or the matrix coordinate ( $M(E)$ -coordinate). We can interpret this inequality as follows:

- the first term  $2W_p(\mu, \nu)$  is to be seen as the initial Euclidean error between the measures  $\mu$  and  $\nu$ ,
- the second term  $\gamma c_{19} \left( \frac{W_p(\mu, \nu)}{r^{d+1}} \right)^{\frac{1}{2}}$  corresponds to the local errors  $W_1(\overline{\mu}_x, \overline{\nu}_y)$  in  $M(E)$  when comparing the normalized covariance matrices of points away from the self-intersections of  $\mathcal{M}$ ,
- the third term  $\gamma c'_{19} \mu(\lambda^r)^{\frac{1}{p}} \left( \frac{W_p(\mu, \nu)}{r^{d+1}} \right)^{\frac{1}{4}}$  stands for the error in  $M(E)$  on points close to the self-intersections of  $\mathcal{M}$ . The quantity of such points is measured via  $\mu(\lambda^r)$ , the measure of the  $r$ -sublevel set of the normal reach.

As a consequence of this proposition, the map  $\mu \mapsto \check{\mu}$ , seen as a map between spaces of measures endowed with the Wasserstein metric, is continuous on the set of measures  $\mu$  which satisfy Hypotheses 5, 6 and 7 with  $\frac{1}{2p} \geq r$ .

We now give the lemmas used in the proof of this Proposition 4.11.

**Lemma 4.12.** *Let  $\pi$  be an optimal transport plan for  $W_p(\mu, \nu)$ . Then*

$$W_{p,\gamma}(\check{\mu}, \check{\nu}) \leq 2^{\frac{p-1}{p}} \left( 1 + \frac{2\gamma}{r} \right) W_p(\mu, \nu) + 2^{\frac{p-1}{p}} \frac{2\gamma}{r} \left( \int W_1^p(\overline{\mu}_x, \overline{\nu}_y) d\pi(x, y) \right)^{\frac{1}{p}}.$$

*Proof.* We first prove the following fact: for every  $x \in \text{supp}(\mu)$  and  $y \in \text{supp}(\nu)$ ,

$$\|\Sigma_\mu(x) - \Sigma_\nu(y)\|_F \leq 2r(\|x - y\| + \mathbf{W}_1(\overline{\mu}_x, \overline{\nu}_y)). \quad (41)$$

Let  $\rho$  be any transport plan between  $\overline{\mu}_x$  and  $\overline{\nu}_y$ . We have

$$\begin{aligned} \Sigma_\mu(x) - \Sigma_\nu(y) &= \int (x - y)^{\otimes 2} d\overline{\mu}_x(x') - \int (y - y')^{\otimes 2} d\overline{\nu}_y(y') \\ &= \int ((x - x')^{\otimes 2} - (y - y')^{\otimes 2}) d\rho(x', y'). \end{aligned} \quad (42)$$

For any  $x' \in \overline{\mathcal{B}}(x, r)$  and  $y' \in \overline{\mathcal{B}}(y, r)$ , we can use Lemma 4.3 to get

$$\left\| (x - x')^{\otimes 2} - (y - y')^{\otimes 2} \right\|_F \leq (r + r)(\|x - y\| + \|x' - y'\|).$$

Therefore, Equation (42) yields

$$\begin{aligned} \|\Sigma_\mu(x) - \Sigma_\nu(y)\|_F &\leq \int 2r(\|x - y\| + \|x' - y'\|) d\rho(x', y') \\ &\leq 2r(\|x - y\| + \mathbf{W}_1(\overline{\mu}_x, \overline{\nu}_y)). \end{aligned}$$

Now, a transport plan  $\pi$  for  $\mathbf{W}_p(\mu, \nu)$  begin given, we build a transport plan  $\tilde{\pi}$  for  $(\check{\mu}, \check{\nu})$  as follows: for every  $\phi: (E \times \mathbf{M}(E))^2 \rightarrow \mathbb{R}$  with compact support, let  $\tilde{\pi}$  satisfy

$$\int \phi(x, A, y, B) d\tilde{\pi}(x, A, y, B) = \int \phi(x, \overline{\Sigma}_\mu(x), y, \overline{\Sigma}_\nu(y)) d\pi(x, y).$$

We have the upper bound

$$\begin{aligned} \mathbf{W}_{p, \gamma}^p(\check{\mu}, \check{\nu}) &\leq \int \|(x, A) - (y, B)\|_\gamma^p d\tilde{\pi}(x, A, y, B) \\ &= \int \left( \|x - y\|^2 + \gamma^2 \|\overline{\Sigma}_\mu(x) - \overline{\Sigma}_\nu(y)\|_F^2 \right)^{\frac{p}{2}} d\pi(x, y) \\ &\leq \int (\|x - y\| + \gamma \|\overline{\Sigma}_\mu(x) - \overline{\Sigma}_\nu(y)\|_F)^p d\pi(x, y) \end{aligned} \quad (43)$$

Besides, Equation (41) gives

$$\|\overline{\Sigma}_\mu(x) - \overline{\Sigma}_\nu(y)\|_F \leq \frac{1}{r^2} \|\Sigma_\mu(x) - \Sigma_\nu(y)\|_F \leq \frac{2}{r} (\|x - y\| + \mathbf{W}_1(\overline{\mu}_x, \overline{\nu}_y)).$$

We can use the inequality  $(a + b)^p \leq 2^{p-1}(a^p + b^p)$ , where  $a, b \geq 0$ , to deduce

$$\begin{aligned} (\|x - y\| + \gamma \|\overline{\Sigma}_\mu(x) - \overline{\Sigma}_\nu(y)\|_F)^p &\leq \left( \|x - y\| + \gamma \frac{2}{r} (\|x - y\| + \mathbf{W}_1(\overline{\mu}_x, \overline{\nu}_y)) \right)^p \\ &\leq 2^{p-1} \left( \left( 1 + \frac{2\gamma}{r} \right) \|x - y\| \right)^p + 2^{p-1} \left( \frac{2\gamma}{r} \mathbf{W}_1(\overline{\mu}_x, \overline{\nu}_y) \right)^p. \end{aligned}$$

By inserting this inequality in Equation (43) we obtain

$$\begin{aligned} \mathbf{W}_{p, \gamma}^p(\check{\mu}, \check{\nu}) &\leq 2^{p-1} \int \left( \left( 1 + \frac{2\gamma}{r} \right) \|x - y\| \right)^p + \left( \frac{2\gamma}{r} \mathbf{W}_1(\overline{\mu}_x, \overline{\nu}_y) \right)^p d\pi(x, y) \\ &= 2^{p-1} \left( 1 + \frac{2\gamma}{r} \right)^p \mathbf{W}_p^p(\mu, \nu) + 2^{p-1} \left( \frac{2\gamma}{r} \right)^p \int \mathbf{W}_1^p(\overline{\mu}_x, \overline{\nu}_y) d\pi(x, y), \end{aligned}$$

which yields the result.  $\square$

**Lemma 4.13.** Let  $w = W_p(\mu, \nu)$  and define  $\alpha = (\frac{w}{r^{d-1}})^{\frac{1}{2}}$ . Suppose that  $r \leq \frac{1}{2p}$  and  $w \leq \min(c_9, 1) (\frac{r}{4})^{d+1}$ . Let  $\pi$  be an optimal transport plan for  $W_p(\mu, \nu)$ . Then

$$\begin{aligned} & \left( \int W_1^p(\overline{\mu}_x, \overline{\nu}_y) d\pi(x, y) \right)^{\frac{1}{p}} \\ & \leq 2^{\frac{p-1}{p}} \left( c_{18} r^{\frac{1}{2}} \mu(\lambda^r)^{\frac{1}{p}} \alpha^{\frac{1}{2}} + \left( 2r^d + c_{16} r^{\frac{d+1}{2}} + c_{17} \right) \alpha + (1 + c_{15})w \right). \end{aligned}$$

If we suppose that  $r \leq 1$ , then

$$\left( \int W_1^p(\overline{\mu}_x, \overline{\nu}_y) d\pi(x, y) \right)^{\frac{1}{p}} \leq 2^{\frac{p-1}{p}} \left( c_{18} r^{\frac{1}{2}} \mu(\lambda^r)^{\frac{1}{p}} \alpha^{\frac{1}{2}} + c_{20} \alpha \right)$$

with  $c_{20} = 3 + c_{15} + c_{16} + c_{17}$ .

*Proof.* We denote  $w = W_p(\mu, \nu)$  and  $\alpha = (\frac{w}{r^{d-1}})^{\frac{1}{2}}$ . Let us subdivide the integral as follows:

$$\int W_1^p(\overline{\mu}_x, \overline{\nu}_y) d\pi(x, y) = \int_A + \int_B + \int_C W_1^p(\overline{\mu}_x, \overline{\nu}_y) d\pi(x, y) \quad (44)$$

where

- $A = \{(x, y) \mid \|x - y\| \geq \alpha\}$ ,
- $B = \{(x, y) \mid \|x - y\| < \alpha \text{ and } \lambda(x) > r\}$ ,
- $C = \{(x, y) \mid \|x - y\| < \alpha \text{ and } \lambda(x) \leq r\}$ .

*Term A.* We use the following simple upper bound:

$$\begin{aligned} W_1(\overline{\mu}_x, \overline{\nu}_y) & \leq W_1(\overline{\mu}_x, \delta_x) + W_1(\delta_x, \delta_y) + W_1(\delta_y, \overline{\nu}_y) \\ & \leq r + \|x - y\| + r \end{aligned}$$

to obtain  $W_1^p(\overline{\mu}_x, \overline{\nu}_y) \leq 2^{p-1} ((2r)^p + \|x - y\|^p)$  and

$$\begin{aligned} \int_A W_1^p(\overline{\mu}_x, \overline{\nu}_y) d\pi(x, y) & \leq \int_A 2^{p-1} ((2r)^p + \|x - y\|^p) d\pi(x, y) \\ & \leq 2^{p-1} (2r)^p \pi(A) + \int 2^{p-1} \|x - y\|^p d\pi(x, y) \\ & = 2^{p-1} (2r)^p \pi(A) + 2^{p-1} w^p. \end{aligned}$$

Besides, Markov inequality yields

$$\pi(A) = \pi(\{(x, y) \mid \|x - y\| > \alpha\}) = \pi(\{(x, y) \mid \|x - y\|^p > \alpha^p\}) \leq \left( \frac{w}{\alpha} \right)^p.$$

Therefore,

$$\begin{aligned} \int_A W_1^p(\overline{\mu}_x, \overline{\nu}_y) d\pi(x, y) & \leq 2^{p-1} (2r)^p \left( \frac{w}{\alpha} \right)^p + 2^{p-1} w^p \\ & = 2^{p-1} (2r^d \alpha)^p + 2^{p-1} w^p, \end{aligned}$$

where we used  $r \frac{w}{\alpha} = r^d \alpha$  on the last line.

*Term B.* On the event  $B$ , we write

$$W_1(\overline{\mu}_x, \overline{\nu}_y) \leq W_1(\overline{\mu}_x, \overline{\mu}_y) + W_1(\overline{\mu}_y, \overline{\nu}_y).$$

Since  $\lambda(x) > r$ , Lemma 4.5 and Lemma 4.7 give  $W_1(\overline{\mu_x}, \overline{\mu_y}) \leq c_{15} \|x - y\|$  and  $W_1(\overline{\mu_y}, \overline{v_y}) \leq c_{17} \alpha$ . We deduce that

$$\begin{aligned} \int_B W_1^p(\overline{\mu_x}, \overline{v_y}) d\pi(x, y) &\leq 2^{p-1} \int_B (c_{15} \|x - y\|)^p + (c_{17} \alpha)^p d\pi(x, y) \\ &\leq 2^{p-1} (c_{15} w)^p + 2^{p-1} (c_{17} \alpha)^p. \end{aligned}$$

*Term C.* We proceed as for Term B, but using Lemmas 4.6 and 4.8 instead of Lemmas 4.5 and 4.7. This yields

$$\begin{aligned} W_1(\overline{\mu_x}, \overline{v_y}) &\leq W_1(\overline{\mu_x}, \overline{\mu_y}) + W_1(\overline{\mu_y}, \overline{v_y}) \\ &\leq c_{16} r^{\frac{1}{2}} \|x - y\|^{\frac{1}{2}} + c_{18} r^{\frac{1}{2}} \alpha^{\frac{1}{2}}, \end{aligned}$$

and we deduce that

$$\begin{aligned} \int_C W_1^p(\overline{\mu_x}, \overline{v_y}) d\pi(x, y) &\leq \int_C 2^{p-1} \left( c_{16} r^{\frac{1}{2}} \|x - y\|^{\frac{1}{2}} \right)^p d\pi(x, y) \\ &\quad + 2^{p-1} \pi(C) \left( c_{18} r^{\frac{1}{2}} \alpha^{\frac{1}{2}} \right)^p. \end{aligned} \tag{45}$$

On the one hand, we have  $\int_C \|x - y\|^{\frac{p}{2}} d\pi(x, y) \leq \int_{E \times E} \|x - y\|^{\frac{p}{2}} d\pi(x, y)$ , and by Jensen's inequality,

$$\int_{E \times E} \|x - y\|^{\frac{p}{2}} d\pi(x, y) \leq (w^p)^{\frac{1}{2}}.$$

On the other hand, by definition of  $C$ , we have  $\pi(C) \leq \mu(\lambda^r)$ . Hence Equation (45) yields

$$\int_C W_1^p(\overline{\mu_x}, \overline{v_y}) d\pi(x, y) \leq 2^{p-1} \left( c_{16} r^{\frac{1}{2}} w^{\frac{1}{2}} \right)^p + 2^{p-1} \mu(\lambda^r) \left( c_{18} r^{\frac{1}{2}} \alpha^{\frac{1}{2}} \right)^p.$$

To conclude the proof, we write

$$\begin{aligned} \int W_1(\overline{\mu_x}, \overline{v_y}) d\pi(x, y) &= \int_A + \int_B + \int_C W_1(\overline{\mu_x}, \overline{v_y}) d\pi(x, y) \\ &\leq 2^{p-1} (2r^d \alpha)^p + 2^{p-1} w^p + 2^{p-1} (c_{15} w)^p + 2^{p-1} (c_{17} \alpha)^p \\ &\quad + 2^{p-1} \left( c_{16} r^{\frac{1}{2}} w^{\frac{1}{2}} \right)^p + 2^{p-1} \mu(\lambda^r) \left( c_{18} r^{\frac{1}{2}} \alpha^{\frac{1}{2}} \right)^p. \end{aligned}$$

We use the inequality  $(a + b)^{\frac{1}{p}} \leq a^{\frac{1}{p}} + b^{\frac{1}{p}}$ , where  $a, b \geq 0$ , to get

$$\begin{aligned} &\left( \int W_1(\overline{\mu_x}, \overline{v_y}) d\pi(x, y) \right)^{\frac{1}{p}} \\ &\leq 2^{\frac{p-1}{p}} \left( 2r^d \alpha + w + c_{15} w + c_{17} \alpha + c_{16} r^{\frac{1}{2}} w^{\frac{1}{2}} + \mu(\lambda^r)^{\frac{1}{p}} c_{18} r^{\frac{1}{2}} \alpha^{\frac{1}{2}} \right) \\ &\leq 2^{\frac{p-1}{p}} \left( c_{18} r^{\frac{1}{2}} \mu(\lambda^r)^{\frac{1}{p}} \alpha^{\frac{1}{2}} + \left( 2r^d + c_{16} r^{\frac{d+1}{2}} + c_{17} \right) \alpha + (1 + c_{15}) w \right), \end{aligned}$$

where we used  $c_{16} r^{\frac{1}{2}} w^{\frac{1}{2}} = c_{16} r^{\frac{d+1}{2}} \alpha$  on the last line. This proves the first result.

If we suppose  $r \leq 1$ , we can use the inequalities  $r^d \leq r^{\frac{d+1}{2}} \leq 1$  and  $w = \alpha r^{\frac{d-1}{2}} w^{\frac{1}{2}} \leq \alpha$  to obtain the simplified expression

$$\left( \int W_1(\overline{\mu_x}, \overline{v_y}) d\pi(x, y) \right)^{\frac{1}{p}} \leq 2^{\frac{p-1}{p}} \left( c_{18} r^{\frac{1}{2}} \mu(\lambda^r)^{\frac{1}{p}} \alpha^{\frac{1}{2}} + (3 + c_{15} + c_{16} + c_{17}) \alpha \right),$$

as wanted. □

## 4.5 An approximation theorem

We are now able to state that the lifted measure  $\check{\nu}$  is close to the exact lifted measure  $\check{\mu}_0$ , that is,  $\check{\nu}$  is a consistent estimator of  $\check{\mu}_0$ , in Wasserstein distance.

**Theorem 4.14.** *Assume that  $\mathcal{M}_0$  and  $\mu_0$  satisfy Hypotheses 1, 2, 3. Let  $\nu$  be any probability measure. Denote  $w = W_p(\mu, \nu)$ . Suppose that  $r \leq \min\left(\frac{1}{2p}, 1\right)$  and  $w \leq \min(c_9, 1)\left(\frac{r}{4}\right)^{d+1}$ . Then*

$$W_{p,\gamma}(\check{\nu}, \check{\mu}_0) \leq \gamma c_{21} \mu(\lambda^r)^{\frac{1}{p}} + \gamma c_{14} r + \gamma c_{19} \left(\frac{w}{r^{d+1}}\right)^{\frac{1}{2}} + 2w$$

where  $c_{21} = 2 + \frac{1}{2}c'_{19}$ .

*Proof.* By using the triangle inequality  $W_{p,\gamma}(\check{\nu}, \check{\mu}_0) \leq W_{p,\gamma}(\check{\nu}, \check{\mu}) + W_{p,\gamma}(\check{\mu}, \check{\mu}_0)$ , we see that the result is a direct consequence of Propositions 4.2 and 4.11.  $\square$

*Remark 4.15.* The quality of the bound given by Theorem 4.14 is balanced by the contributions of Propositions 4.2 and 4.11. According to the first one, the quantity  $W_{p,\gamma}(\check{\mu}, \check{\mu}_0)$  is minimized when  $r$  is as small as possible, and according to the second one, the distance  $W_{p,\gamma}(\check{\nu}, \check{\mu})$  is minimized when  $r$  is chosen large. Roughly speaking, to optimize the bound given by the theorem, we have to pick a  $r$  given by equating the terms  $r$  and  $\left(\frac{w}{r^{d+1}}\right)^{\frac{1}{2}}$ , that is,  $r = w^{\frac{1}{d+3}}$ . We will make this choice in the following corollary.

*Remark 4.16.* In the case where  $\mathcal{M}_0$  is embedded, we have seen in Proposition 3.6 that the normal reach  $\lambda$  is bounded below by  $\text{reach}(\mathcal{M}) > 0$ . In particular,  $\mu(\lambda^r)$  is zero for  $r$  small enough. In this case, Theorem 4.14 reads

$$W_{p,\gamma}(\check{\nu}, \check{\mu}_0) \leq \gamma c_{14} r + \gamma c_{19} \left(\frac{w}{r^{d+1}}\right)^{\frac{1}{2}} + 2w$$

We deduce an approximation result: if  $(\nu_i)_{i \geq 0}$  is a sequence of probability measures such that  $w_i = W_p(\mu, \nu_i)$  goes to zero, and if we choose a sequence of radii  $(r_i)_{i \geq 0}$  such that  $(r_i)_{i \geq 0}$  and  $(w_i/r_i^{d+1})_{i \geq 0}$  go to zero, then  $W_{p,\gamma}(\check{\nu}_i, \check{\mu}_0)$  goes to zero too.

More generally,  $W_{p,\gamma}(\check{\nu}_i, \check{\mu}_0)$  goes to zero if we only assume that  $\mathcal{M}_0$  satisfies Hypothesis 4. This is stated in the following corollary, which is a weaker version of the theorem, that we shall use in the following section to simplify the results.

**Corollary 4.17.** *Let  $r > 0$ . Assume that  $\mathcal{M}_0$  and  $\mu_0$  satisfy Hypotheses 1, 2 and 4. Let  $\nu$  be any probability measure. Denote  $w = W_p(\mu, \nu)$ . Suppose that  $r < \min\left(\frac{1}{2p}, r_4, 1\right)$  and  $w \leq \min(c_9, 1)\left(\frac{r}{4}\right)^{d+3}$ . Then*

$$W_{p,\gamma}(\check{\nu}, \check{\mu}_0) \leq (1 + \gamma c_{22}) r^{\frac{1}{p}}$$

with  $c_{22} = c_{21}(c_4)^{\frac{1}{p}} + c_{19} + c_{14}$ .

*Proof.* According to Theorem 4.14, we have

$$W_{p,\gamma}(\check{\nu}, \check{\mu}_0) \leq \gamma c_{21} \mu(\lambda^r)^{\frac{1}{p}} + \gamma c_{14} r + \gamma c_{19} \left(\frac{w}{r^{d+1}}\right)^{\frac{1}{2}} + 2w.$$

Note that the assumption  $w \leq \left(\frac{r}{4}\right)^{d+3}$  yields  $\left(\frac{w}{r^{d+1}}\right)^{\frac{1}{2}} \leq r$ . Besides,  $r \leq 1$  yields  $w \leq \left(\frac{r}{4}\right)^{d+3} \leq \frac{r}{2}$ . Finally, Hypothesis 4 gives  $\mu(\lambda^r) \leq c_4 r$ , and we deduce the result:

$$\begin{aligned} W_{p,\gamma}(\check{\nu}, \check{\mu}_0) &\leq \gamma c_{21} (c_4 r)^{\frac{1}{p}} + \gamma c_{14} r + \gamma c_{19} r + r \\ &\leq \left(\gamma c_{21} (c_4)^{\frac{1}{p}} + \gamma c_{14} + \gamma c_{19} + 1\right) r^{\frac{1}{p}} \end{aligned}$$

where we used to the weak upper bound  $r \leq r^{\frac{1}{p}}$  on the last line.  $\square$

## 5 Topological inference with the lifted measure

Based on the results of the last section, we show how the lifted measure  $\check{\nu}$  can be used to infer the homotopy type of  $\check{\mathcal{M}}$  or its homology groups.

### 5.1 Overview of the method

Let us recall the results obtained so far. Assume that the immersion  $u: \mathcal{M}_0 \rightarrow \mathcal{M}$  and the measure  $\mu_0$  satisfy the Hypotheses 1, 2 and 3. Our goal is to estimate the exact lifted measure  $\check{\mu}_0$  on  $E \times \mathbf{M}(E)$ , since its support is the submanifold  $\check{\mathcal{M}}$ , which is diffeomorphic to  $\mathcal{M}_0$ . To do so, we suppose that we are observing a measure  $\nu$  on  $E$ . No assumptions are made on  $\nu$ . Our results only depends on the Wasserstein distance  $w = W_p(\mu, \nu)$ , where  $\mu = u_*\mu_0$ . Recall that the measure  $\check{\mu}_0$  is defined as (see Equation (21)):

$$\check{\mu}_0 = (u_*\mu_0)(x_0) \otimes \left\{ \delta_{\frac{1}{d+2} p_{T_{x_0}\mathcal{M}}} \right\}.$$

To approximate  $\check{\mu}_0$ , pick a parameter  $r > 0$  and consider the lifted measure  $\check{\nu}$  built on  $\nu$  (see Definition 4.2):

$$\check{\nu} = \nu(x) \otimes \left\{ \delta_{\Sigma_\nu(x)} \right\}.$$

Choose  $\gamma > 0$ . Endow the space  $E \times \mathbf{M}(E)$  with the norm  $\|\cdot\|_\gamma$  (see Equation (22)), and consider the Wasserstein distance  $W_{p,\gamma}(\cdot, \cdot)$  between measures on  $E \times \mathbf{M}(E)$  (see Equation (23)). We quantify the quality of the approximation by the Wasserstein distance  $W_{p,\gamma}(\check{\mu}_0, \check{\nu})$ . According to Theorem 4.14, we have

$$W_{p,\gamma}(\check{\nu}, \check{\mu}_0) \leq \gamma c_{21} \mu(\lambda^r)^{\frac{1}{p}} + \gamma c_{14} r + \gamma c_{19} \left( \frac{w}{r^{d+1}} \right)^{\frac{1}{2}} + 2w$$

as long as the parameter  $r$  satisfies

$$4 \left( \frac{w}{\min(c_9, 1)} \right)^{\frac{1}{d+1}} \leq r \leq \min \left( \frac{1}{2\rho}, 1 \right).$$

Under Hypothesis 4, Corollary 4.17 gives a weaker form of this result. We have

$$W_{p,\gamma}(\check{\nu}, \check{\mu}_0) \leq (1 + \gamma c_{22}) r^{\frac{1}{p}}$$

as long as the parameter  $r$  satisfies

$$4 \left( \frac{w}{\min(c_9, 1)} \right)^{\frac{1}{d+3}} \leq r \leq \min \left( \frac{1}{2\rho}, r_4, 1 \right).$$

In the following subsections, we show how these results lead to consistent estimations of  $\mathcal{M}_0$  and its homology. Namely, we can estimate the homotopy type of  $\check{\mathcal{M}}$ , and hence of  $\mathcal{M}_0$ , by considering the sublevel sets of the DTM  $d_{\check{\nu},m,\gamma}$  (see Corollary 5.3). The notation  $d_{\check{\nu},m,\gamma}$  corresponds to the DTM, as defined in Subsect. 2.3, with measure  $\check{\nu}$ , parameter  $m$ , and seen in the ambient space  $(E \times \mathbf{M}(E), \|\cdot\|_\gamma)$ . Besides, we can estimate the persistent homology of the DTM-filtration  $W_\gamma[\check{\mu}_0]$  with the filtration  $W_\gamma[\check{\nu}]$  (see Corollary 5.5). Here,  $W_\gamma[\cdot]$  corresponds to the DTM-filtration in the ambient space  $(E \times \mathbf{M}(E), \|\cdot\|_\gamma)$ .

**Example 5.1.** Let  $\mathcal{M}$  be the lemniscate of Bernoulli of diameter 2. It is the immersion of a circle  $\mathcal{M}_0$ . We observe a 100-sample  $X$  of  $\mathcal{M}$  (Figure 22). Experimentally, we computed the Hausdorff distance  $d_H(\mathcal{M}, X) \approx 0.026$ . Let  $\mu$  be the Hausdorff measure on  $\mathcal{M}$  and  $\nu$  the empirical measure on  $X$ . We choose the parameter  $p = 2$ . Their Wasserstein distance is approximately  $W_2(\mu, \nu) \approx 0.015$ .



Figure 22: Left: The lemniscate  $\mathcal{M}$ . Right: The set  $X$ , a 100-sample of  $\mathcal{M}$ .

For each point  $x$  of  $X$ , we compute the normalized local covariance matrix  $\bar{\Sigma}_v(x)$  with parameter  $r = 0.5$  and  $0.1$ . This matrix is used as an estimator of the tangent space  $T_x\mathcal{M}$ . In order to observe the quality of this estimation, we represent on Figure 23 (first row) the principal axes of  $\bar{\Sigma}_v(x)$  for some  $x$ . On the second row are represented the distances  $\|\bar{\Sigma}_v(x) - \frac{1}{d+2}P_{T_x\mathcal{M}}\|_F$ . One sees that  $r = 0.1$  yields a better approximation. However, the estimation is still biased next to the self-intersection points of  $\mathcal{M}$ .

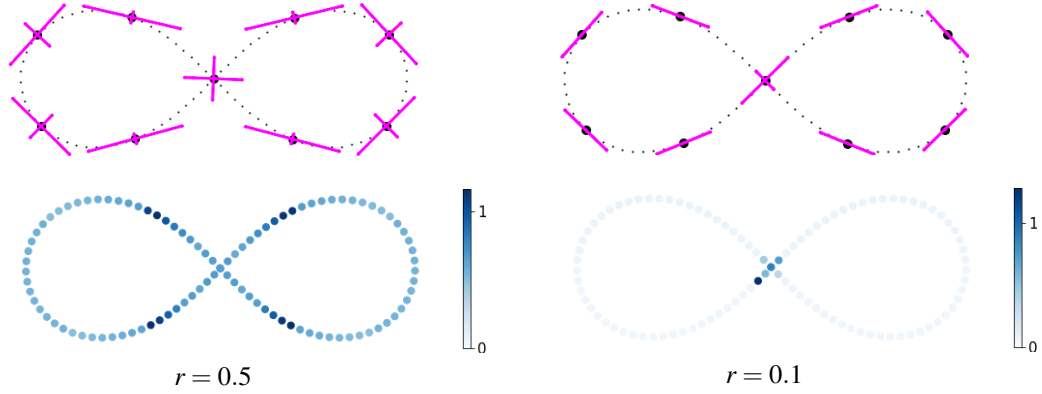


Figure 23: First row: The eigenvectors of  $\bar{\Sigma}_v(x)$  for some  $x \in X$ , weighted with their corresponding eigenvalue. Second row: color representation of the distances  $\|\bar{\Sigma}_v(x) - \frac{1}{d+2}P_{T_x\mathcal{M}}\|_F$ .

Now we choose the parameter  $\gamma = 2$ . For  $r = 0.5$  and  $0.1$ , we consider the lifted measures built on  $v$ , respectively denoted  $\check{v}^{0.5}$  and  $\check{v}^{0.1}$ . They are measure on the lift space  $\mathbb{R}^2 \times \mathbf{M}(\mathbb{R}^2)$ , which is endowed with the norm  $\|\cdot\|_\gamma$ . We computed the Wasserstein distances:

$$W_{2,\gamma}(\check{\mu}_0, \check{v}^{0.5}) \approx 0.674 \quad \text{and} \quad W_{2,\gamma}(\check{\mu}_0, \check{v}^{0.1}) \approx 0.200.$$

In comparison, even with a small parameter  $r$ , the Hausdorff distance between their support is still large:

$$d_H(\check{\mathcal{M}}, \text{supp}(\check{v}^{0.5})) \approx 1.142 \quad \text{and} \quad d_H(\check{\mathcal{M}}, \text{supp}(\check{v}^{0.1})) \approx 1.273.$$

These sets are represented in Figure 24. Observe that, at the center of the graphs, the measures  $\check{v}^{0.5}$  and  $\check{v}^{0.1}$  deviate from the set  $\check{\mathcal{M}}$ .

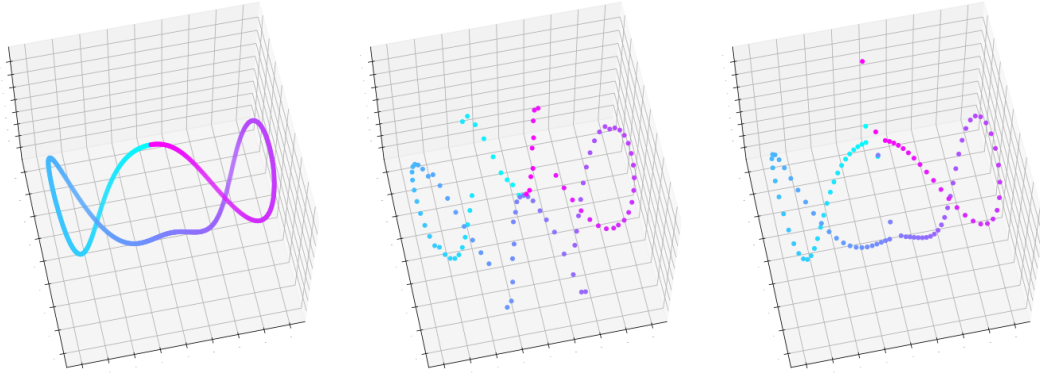


Figure 24: Left: The lifted lemniscate  $\check{\mathcal{M}}$ , projected in a 3-dimensional subspace via PCA. Center: The set  $\text{supp}(\check{\nu}^{0.5})$  projected in the same 3-dimensional subspace. Right: Same for  $\text{supp}(\check{\nu}^{0.1})$ .

**Example 5.2.** Let  $u: \mathcal{M}_0 \rightarrow \mathcal{M}$  be the figure-8 immersion of the torus in  $\mathbb{R}^3$ , represented in Figure 25. It can be parametrized by rotating a lemniscate around an axis, while forming a full twist. The self-intersection points of this immersion corresponds to the inner circle formed by the center of the lemniscate. These are the points  $x$  of  $\mathcal{M}$  such that their normal reach  $\lambda(x)$  is zero.

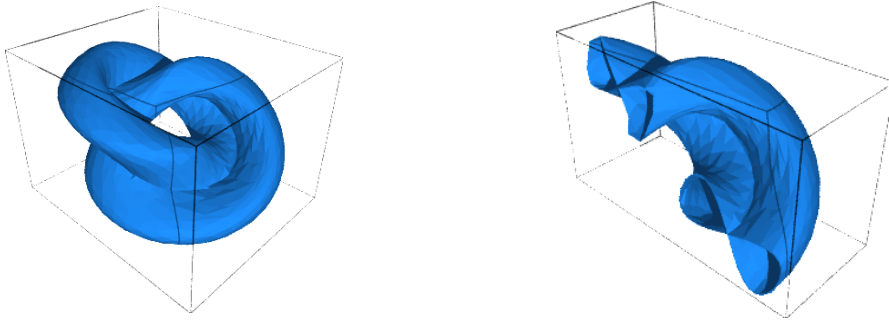


Figure 25: Left: The immersion  $\mathcal{M}$  of the torus. Right: A section of  $\mathcal{M}$ . One sees the inner lemniscate.

Let  $\check{\mathcal{M}}$  be the lift of  $\mathcal{M}_0$ . It is a submanifold of  $\mathbb{R}^3 \times \mathbf{M}(\mathbb{R}^3) \simeq \mathbb{R}^{12}$ . One cannot embed  $\check{\mathcal{M}}$  in  $\mathbb{R}^3$  by performing a PCA. However, we can try to visualize  $\check{\mathcal{M}}$  by considering a small section of it. Figure 26 represents a subset of  $\check{\mathcal{M}}$ , projected in a 3-dimensional subspace via PCA. One sees that it does not self-intersect.

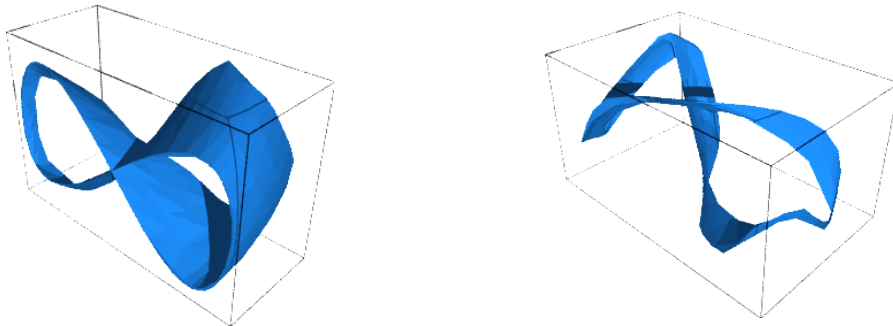


Figure 26: Left: A section of  $\mathcal{M}$ . Right: The corresponding section of  $\check{\mathcal{M}}$ , projected in a 3-dimensional subspace via PCA. Observe that it does not self-intersect.

In order to fit in the context of our study, let  $\mu$  be the Hausdorff measure on  $\mathcal{M}$ . We observe a 9000-sample  $X$  of  $\mathcal{M}$ , and consider its empirical measure  $\nu$ . The set  $X$  is depicted in Figure 27. Choose the



parameter  $p = 1$ . We compute the Wasserstein distance  $W_1(\mu, \nu) \approx 0.070$  and the Hausdorff distance  $d_H(\mathcal{M}, X) = 0.083$ . Let  $r = 0.09$ . In order to observe the estimation of tangent spaces by local covariance matrices  $\bar{\Sigma}_\nu(x)$  with parameter  $r$ , we represent on Figure 27 the points  $x$  such that the distance  $\|\bar{\Sigma}_\nu(x) - \frac{1}{d+2}P_{T_x\mathcal{M}}\|_F$  is greater than 1. Observe that the estimation is biased next to the self-intersection circle of  $\mathcal{M}$ .

Last, let us choose the parameter  $\gamma = 2$ , and consider the lifted measure  $\check{\nu}$ . We have  $W_1(\check{\mu}_0, \check{\nu}) \approx 0.986$ . In comparison, the Hausdorff distance between their support is large:  $d_H(\check{\mathcal{M}}, \text{supp}(\check{\nu})) \approx 2.188$ .

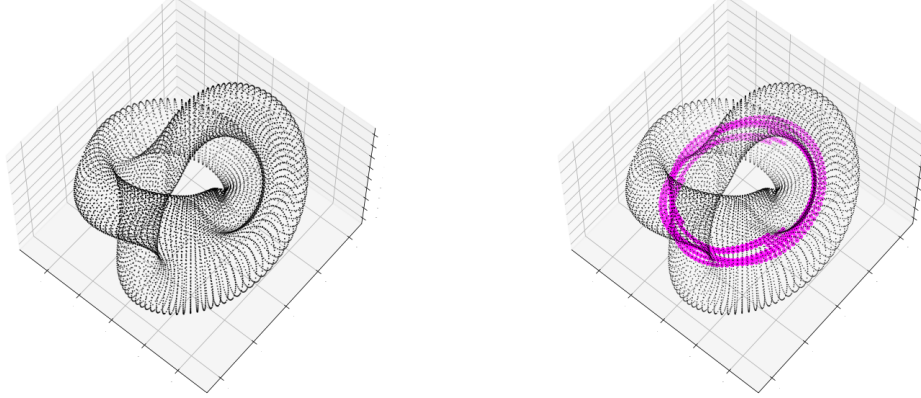


Figure 27: Left: The set  $X$ , a sample of  $\mathcal{M}$ . Right: The set  $X$ , where  $x \in X$  is colored in magenta if  $\|\bar{\Sigma}_\nu(x) - \frac{1}{d+2}P_{T_x\mathcal{M}}\|_F \geq 1$ .

## 5.2 Homotopy type estimation with the DTM

In this subsection, we use the DTM, as defined in Subsect. 2.3, to infer the homotopy type of  $\check{\mathcal{M}}$  from the lifted measure  $\check{\nu}$ . We shall use the DTM on  $\check{\nu}$ , which lives in the space  $E \times M(E)$  endowed with the norm  $\|\cdot\|_\gamma$ . It is denoted  $d_{\check{\nu}, m, \gamma}$ .

In order to apply Theorem 2.7 in our setting, we have to consider geometric quantities associated to the submanifold  $\check{\mathcal{M}}$ . Note that the map  $\check{u}$  itself satisfies the Hypotheses 2 and 3, since the immersion  $u$  does. Hence we can consider the following quantities: for every  $\gamma > 0$ , we denote by

- $\text{reach}_\gamma(\check{\mathcal{M}})$  the reach of  $\check{\mathcal{M}}$  (for the norm  $\|\cdot\|_\gamma$ ),
- $\check{\rho}_\gamma, \check{L}_{0, \gamma}, \check{f}_{\min, \gamma}$  and  $\check{f}_{\max, \gamma}$  the constants given by Hypotheses 2 and 3 applied to  $\check{\mathcal{M}}$ ,
- $\check{c}_{23, \gamma} = \check{f}_{\min, \gamma} J_{\min} V_d$  the constant given by Proposition 3.14 Point 1 applied to  $\check{\mu}_0$ .

According to Subsect. 2.1, a sufficient condition for  $\check{\mathcal{M}}$  to satisfy  $\text{reach}_\gamma(\check{\mathcal{M}}) > 0$  is that it is a  $\mathcal{C}^2$ -submanifold. This would be the case if  $\mathcal{M}_0$  and  $u$  were  $\mathcal{C}^3$ . Also, we point out that the constant  $\check{\rho}_\gamma$  cannot be deduced from  $\rho$ : the first one can be arbitrary large or small compared to the second one, even with  $\gamma$  being fixed. This remark holds for the other constants.

These constants being given, we propose a way to tune the parameters  $r, \gamma, m$  and  $t$  in such a way that the  $t$ -sublevel set  $d_{\check{\nu}, m, \gamma}^t$  of the DTM captures the homotopy type of  $\check{\mathcal{M}}$ , or equivalently, of  $\mathcal{M}_0$ .

**Corollary 5.3.** *Assume that  $\mathcal{M}_0$  and  $\mu_0$  satisfy Hypotheses 1, 2, 3 and 4. Let  $\nu$  be any probability measure on  $E$ . Denote  $w = W_2(\mu, \nu)$ . Choose  $r > 0, \gamma > 0$  and  $m \in (0, 1)$  such that*

- $4 \left( \frac{w}{\min(c_9, 1)} \right)^{\frac{1}{d+3}} \leq r \leq \min \left( \frac{1}{2\rho}, r_4, 1 \right)$
- $m \leq \frac{\check{c}_{23, \gamma}}{(2\check{\rho}_\gamma)^d}$  and

$$\bullet \quad (1 + \gamma c_{22}) r^{\frac{1}{2}} \leq m^{\frac{1}{2}} \left( \frac{\text{reach}_{\gamma}(\check{\mathcal{M}})}{9} - \left( \frac{m}{\check{c}_{23,\gamma}} \right)^{\frac{1}{d}} \right).$$

Define  $\varepsilon$  and choose  $t$  as follows:

$$\varepsilon = \left( \frac{m}{\check{c}_{23,\gamma}} \right)^{\frac{1}{d}} + (1 + \gamma c_{22}) \left( \frac{r}{m} \right)^{\frac{1}{2}} \quad \text{and} \quad t \in [4\varepsilon, \text{reach}_{\gamma}(\check{\mathcal{M}}) - 3\varepsilon].$$

Then the sublevel set of the DTM  $d_{\check{\nu},m,\gamma}^t$  is homotopy equivalent to  $\mathcal{M}_0$ .

*Proof.* In order to fit in the context of Theorem 2.7, we have to consider the usual Euclidean norm  $\|\cdot\|$  on  $E \times M(E)$ . It corresponds to the norm  $\|\cdot\|_{\gamma}$  with  $\gamma = 1$ . For a general parameter  $\gamma > 0$ , consider the dilatation map  $i_{\gamma}: E \times M(E) \rightarrow E \times M(E)$  defined as

$$i_{\gamma}: (x, A) \mapsto (x, \gamma A).$$

A computation shows that, for every probability measures  $\alpha, \beta$  on  $E \times M(E)$ , we have

$$W_{2,\gamma}(\alpha, \beta) = W_2((i_{\gamma})_*\alpha, (i_{\gamma})_*\beta),$$

where  $W_2$  denotes the 2-Wasserstein distance on  $E \times M(E)$  endowed with the usual Euclidean norm  $\|\cdot\|$ . Corollary 4.17 then reads

$$W_2((i_{\gamma})_*\check{\mu}_0, (i_{\gamma})_*\check{\nu}) \leq (1 + \gamma c_{22}) r^{\frac{1}{2}},$$

where  $(i_{\gamma})_*\check{\mu}_0$  and  $(i_{\gamma})_*\check{\nu}$  are the push-forwards of  $\check{\mu}_0$  and  $\check{\nu}$  by the map  $i_{\gamma}$ . Besides, consider the set

$$\check{\mathcal{M}}_{\gamma} = i_{\gamma}(\check{\mathcal{M}}) = \{(x, \gamma A) \mid (x, A) \in \check{\mathcal{M}}\}.$$

It is clear that

$$\text{reach}_{\gamma}(\check{\mathcal{M}}) = \text{reach}(\check{\mathcal{M}}_{\gamma}),$$

where we recall that  $\text{reach}_{\gamma}(\check{\mathcal{M}})$  is the reach of  $\check{\mathcal{M}}$  with respect to the norm  $\|\cdot\|_{\gamma}$ , and  $\text{reach}(\check{\mathcal{M}}_{\gamma})$  is the reach of  $\check{\mathcal{M}}_{\gamma}$  with respect to the usual norm  $\|\cdot\|$  on  $E \times M(E)$ . Finally, consider the DTM  $d_{(i_{\gamma})_*\check{\nu},m}$  with respect to the usual Euclidean norm. Observe that, for every  $t \geq 0$ , the sublevel sets of the DTM  $d_{(i_{\gamma})_*\check{\nu},m}$  and  $d_{\check{\nu},m,\gamma}^t$  are linked via

$$d_{\check{\nu},m,\gamma}^t = i_{\gamma} \left( d_{\check{\nu},m,\gamma}^t \right).$$

In particular, they share the same homotopy type. Now we obtain the result as a consequence of Theorem 2.7 applied to the measures  $(i_{\gamma})_*\check{\mu}_0$  and  $(i_{\gamma})_*\check{\nu}$ . Let us verify that the assumptions of the theorem are satisfied. Our assumption about  $m$  ensures that

$$\left( \frac{m}{\check{c}_{23,\gamma}} \right)^{\frac{1}{d}} \leq \frac{1}{2\check{\rho}_{\gamma}},$$

hence by Proposition 3.14 Point 1 we get  $\check{\mu}_0(\mathcal{B}(x, r)) \geq \check{c}_{23,\gamma} r^d$  for all  $x \in \text{supp}(\check{\mu}_0)$  and  $r < \left( \frac{m}{\check{c}_{23,\gamma}} \right)^{\frac{1}{d}}$ .

Moreover, the assumption about  $(1 + \gamma c_{22}) r^{\frac{1}{2}}$  ensures that

$$W_2((i_{\gamma})_*\check{\mu}_0, (i_{\gamma})_*\check{\nu}) \leq m^{\frac{1}{2}} \left( \frac{\text{reach}_{\gamma}(\check{\mathcal{M}})}{9} - \left( \frac{m}{\check{c}_{23,\gamma}} \right)^{\frac{1}{d}} \right)$$

is satisfied, since  $W_2((i_{\gamma})_*\check{\mu}_0, (i_{\gamma})_*\check{\nu}) \leq (1 + \gamma c_{22}) r^{\frac{1}{2}}$  by Corollary 4.17.  $\square$

**Example 5.4.** Let  $\mathcal{M}$  be the lemniscate of Bernoulli, as in Example 5.1. Suppose that  $\mu$  is the uniform distribution on  $\mathcal{M}$ , and  $\nu$  is the empirical measure on a 500-sample of  $\mathcal{M}$ . We choose the parameters  $\gamma = 2$ ,  $r = 0.03$  and  $m = 0.01$ . Let  $\check{\nu}$  be the lifted measure associated to  $\nu$ . Figure 28 represents set the  $\text{supp}(\check{\nu})$ , and the values of the DTM  $d_{\check{\nu},m,\gamma}$  on it. Observe that the anomalous points, i.e., points for which the local covariance matrix is not well estimated, have large DTM values.

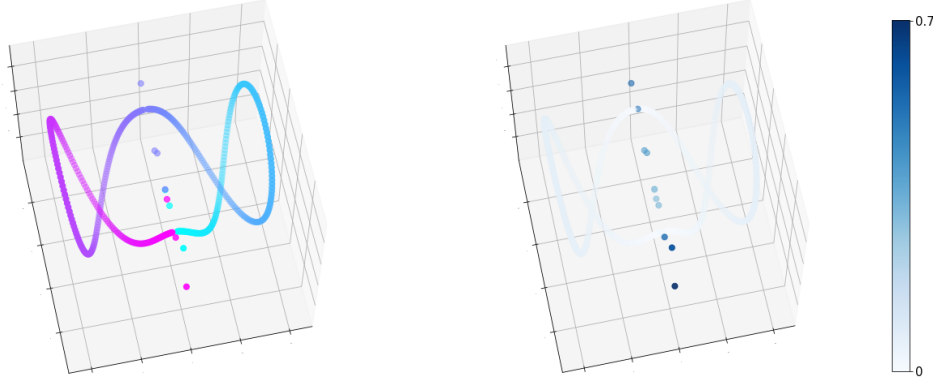


Figure 28: Left: The set  $\text{supp}(\check{\nu}) \subset \mathbb{R}^2 \times \mathbf{M}(\mathbb{R}^2)$ , projected in a 3-dimensional subspace via PCA. Right: The set  $\text{supp}(\check{\nu})$  with colors indicating the value of the DTM  $d_{\check{\nu},m,\gamma}$ .

### 5.3 Persistent homology with DTM-filtrations

In this subsection, we aim to estimate the DTM-filtration of  $\check{\mu}_0$ , as defined in Subsect. 2.3, from  $\nu$ . We shall use the DTM-filtration on  $\check{\nu}$ , denoted  $W_\gamma[\check{\nu}]$ , with respect to the ambient norm  $\|\cdot\|_\gamma$  on  $E \times \mathbf{M}(E)$ . We use the notations  $\check{\rho}_\gamma$  and  $\check{c}_{23,\gamma}$  of the previous subsection.

**Corollary 5.5.** Assume that  $\mathcal{M}_0$  and  $\mu_0$  satisfy Hypotheses 1, 2, 3 and 4. Let  $\nu$  be any probability measure. Denote  $W_2(\mu, \nu) = w$ . Choose  $r > 0$ ,  $\gamma > 0$  and  $m \in (0, 1)$  such that

- $4 \left( \frac{w}{\min(c_9, 1)} \right)^{\frac{1}{d+3}} \leq r \leq \min \left( \frac{1}{2p}, r_4, 1 \right),$
- $m \leq \frac{\check{c}_{23,\gamma}}{(2\check{\rho}_\gamma)^d},$
- $(1 + \gamma c_{22}) r^{\frac{1}{2}} \leq \frac{1}{4}.$

Then we have a bound on the interleaving distance between the DTM-filtrations:

$$d_i(W_\gamma[\check{\mu}_0], W_\gamma[\check{\nu}]) \leq \check{c}_{1,\gamma} (1 + \gamma c_{22})^{\frac{1}{2}} m^{-\frac{1}{2}} r^{\frac{1}{4}} + \check{c}'_{1,\gamma} m^{\frac{1}{d}},$$

where  $\check{c}_{1,\gamma} = 8 \text{diam}(\mathcal{M}) + 8\gamma + 5$  and  $\check{c}'_{1,\gamma} = 2 (\check{c}_{23,\gamma})^{-\frac{1}{d}}$ .

*Proof.* As in the proof of Corollary 5.3, let  $i_\gamma$  be the map  $i_\gamma: (x, A) \mapsto (x, \gamma A)$ . Let  $W[\cdot]$  denotes the DTM-filtration on  $\check{\nu}$  with respect to the usual Euclidean norm. That is, the filtration  $W[\cdot]$  corresponds to  $W_\gamma[\cdot]$  with  $\gamma = 1$ . A computation shows that the filtration  $W[(i_\gamma)_* \check{\nu}]$  and  $W_\gamma[\check{\nu}]$  are linked via

$$W[(i_\gamma)_* \check{\nu}] = i_\gamma(W_\gamma[\check{\nu}]).$$

Now let  $\check{w} = W_2((i^\gamma)_* \check{\mu}_0, (i^\gamma)_* \check{\nu})$ . We have  $\check{w} = W_{2,\gamma}(\check{\mu}_0, \check{\nu})$ , hence Corollary 4.17 gives

$$\check{w} \leq (1 + \gamma c_{22}) r^{\frac{1}{2}}. \quad (46)$$

Moreover, we can apply Corollary 2.9 to  $\mu = (i^\gamma)_* \check{\mu}_0$  and  $\nu = (i^\gamma)_* \check{\nu}$  to get

$$d_i(W[(i^\gamma)_* \check{\mu}_0], W[(i^\gamma)_* \check{\nu}]) \leq \check{c}_{1,\gamma} \left( \frac{\check{w}}{m} \right)^{\frac{1}{2}} + \check{c}'_{1,\gamma} m^{\frac{1}{d}}, \quad (47)$$

where  $\check{c}_{1,\gamma} = (8 \text{diam}(\check{\mathcal{M}}) + 5)$  and  $\check{c}'_{1,\gamma} = 2 (\check{c}_{23,\gamma})^{-\frac{1}{d}}$ . Note that

$$\text{diam}(\check{\mathcal{M}}) \leq \left( \text{diam}(\mathcal{M})^2 + \gamma^2 \left( 2 \frac{1}{2} \right)^2 \right)^{\frac{1}{2}} \leq \text{diam}(\mathcal{M}) + \gamma$$

since the matrices  $\frac{1}{d+2} p_{T_x, \mathcal{M}}$  have norm  $\left\| \frac{1}{d+2} p_{T_x, \mathcal{M}} \right\|_F = \frac{\sqrt{d}}{d+2} \leq \frac{1}{2}$ . Our assumption  $m \leq \frac{\check{c}_{23,\gamma}}{(2\check{\rho}_\gamma)^d}$  ensures that the condition  $\check{\mu}_0(\mathcal{B}(x, r)) \geq \check{c}_{23,\gamma} r^d$  of the theorem is satisfied. Similarly, the assumption  $(1 + \gamma c_{22}) r^{\frac{1}{2}} \leq \frac{1}{4}$  yields  $\check{w} \leq \frac{1}{4}$ .

Combining Equations (46) and (47) we get

$$d_i(W[(i^\gamma)_* \check{\mu}_0], W[(i^\gamma)_* \check{\nu}]) \leq \check{c}_{1,\gamma} (1 + \gamma c_{22})^{\frac{1}{2}} m^{-\frac{1}{2}} r^{\frac{1}{4}} + \check{c}'_{1,\gamma} m^{\frac{1}{d}}.$$

Now, by using the definition of an interleaving of filtrations, one proves that

$$d_i(W_\gamma[\check{\mu}_0], W_\gamma[\check{\nu}]) = d_i(W[(i^\gamma)_* \check{\mu}_0], W[(i^\gamma)_* \check{\nu}]),$$

and we obtain the result.  $\square$

**Example 5.6.** Say that  $\mu$  is the uniform measure on the union of five intersecting circles of radius 1. We observe  $\nu$ , the empirical measure on the point cloud  $X$  drawn in Figure 29. It consists of 300 points per circle, and 100 anomalous points. Let  $p = 1$ . Experimentally, we have  $W_1(\mu, \nu) \approx 0.044$ .

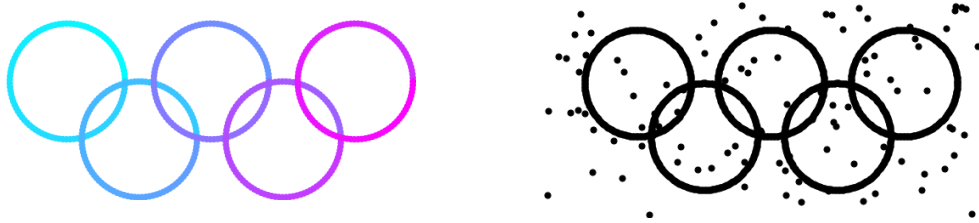


Figure 29: Left: the set  $\mathcal{M} = \text{supp}(\mu)$ . Right: The set  $X = \text{supp}(\nu)$ .

Let  $\gamma = 1$ . Observe that the barcodes of the DTM-filtration of the exact lifted measure  $W_\gamma[\check{\mu}_0]$ , represented in Figure 30, reveal the homology of the disjoint union of five circles—which is the set  $\mathcal{M}_0$ . Only bars of length larger than 0.1 are displayed. We consider the construction of the lifted  $\check{\nu}$  with parameter  $r = 0.03$ , and the DTM-filtration with  $m = 0.01$ . The barcodes of the DTM-filtration  $W_\gamma[\check{\nu}]$  are close to the barcodes of  $W_\gamma[\check{\mu}_0]$ . To compare, we also plot the persistence barcodes of the usual Čech filtration on  $\text{supp}(\check{\nu})$ . Observe that the five connected components do not appear clearly anymore.

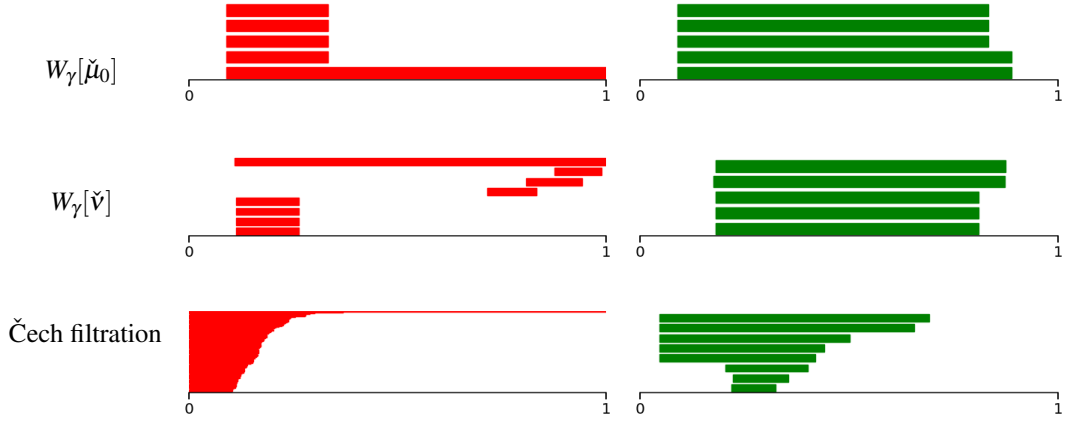


Figure 30: First row: Persistence barcode of the 0- and 1-homology of the DTM-filtration on  $\check{\mu}_0$ . Second row: Same for  $\check{\nu}$ . Third row: Persistence barcodes of the usual Čech filtration on  $\text{supp}(\check{\nu})$ .

At this point, we can propose a clustering procedure based on  $W_\gamma[\check{\nu}]$ . First, select a  $t \in [0, +\infty)$ . Then, extract the connected components of the set  $W_\gamma^t[\check{\nu}]$  of the DTM-filtration. We show in Figure 31 the components we obtain for several values of  $t$ . We see that there exists a value for which the five circles are well clustered ( $t = 0.2$ ). Besides, observe that small values of  $t$  (resp. large) may lead to more connected components than wanted (resp. less).

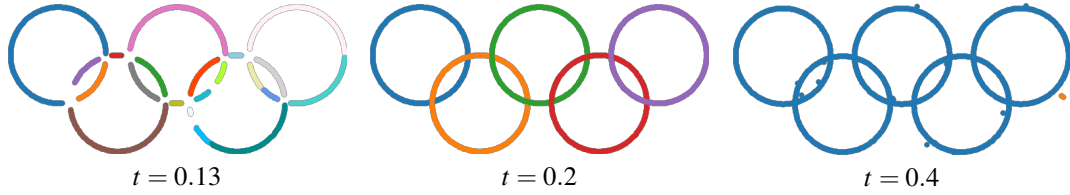


Figure 31: Components obtained by the clustering procedure, where each color correspond to a cluster. The clusterings consist respectively in 21, 5 and 2 clusters.

From an algorithmic viewpoint, this clustering can be obtained by computing the connected components of the nerve of the set  $W_\gamma^t[\check{\nu}]$  or, equivalently, the connected components of its underlying graph  $G$ . As we see from the definition of the DTM-filtration (Equation (4)), the vertices of  $G$  are the points  $\check{x} \in \text{supp}(\check{\nu})$  with DTM value  $d_{\check{\nu},m,\gamma}(\check{x})$  not greater than  $t$ , and where an edge  $[\check{x}, \check{y}]$  is added if

$$\|x - y\| + d_{\check{\nu},m,\gamma}(\check{x}) + d_{\check{\nu},m,\gamma}(\check{y}) \leq t.$$

**Example 5.7.** Consider the immersion of the Klein bottle in  $\mathbb{R}^3$  represented in Figure 32. Note that the self-intersection of this immersion forms a circle. We consider a  $22^1092$ -sample  $X$  of it.

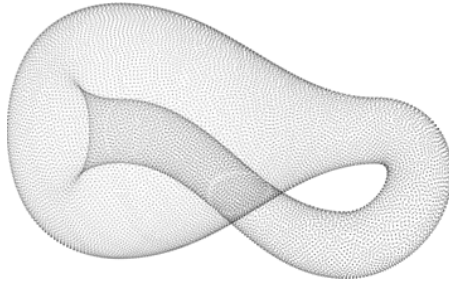


Figure 32: Sample of the Klein bottle immersed in  $\mathbb{R}^3$ .

Let  $\nu$  be the empirical measure on this point cloud. We build the lifted measure  $\check{\nu}$  with parameters  $r = 0.08$  and  $\gamma = 3$ , and we consider the DTM-filtration  $W_\gamma[\check{\nu}]$  with parameter  $m = 0.0001$ . The barcodes of this filtration are depicted in Figure 33, with coefficients in two finite field:  $\mathbb{Z}/2\mathbb{Z}$  and  $\mathbb{Z}/3\mathbb{Z}$ . We also plot the barcodes of the usual Čech filtration of  $X$  in  $\mathbb{R}^3$ . Only bars of length larger than 0.4 are displayed.

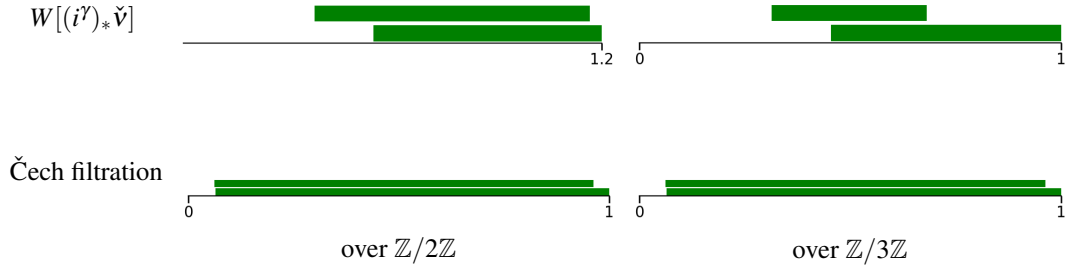


Figure 33: First row: Persistence barcode of the 1-homology of the DTM-filtration on  $\check{\nu}$ . Second row: Persistence barcodes of the usual Čech filtration on  $X$ .

We see that the barcodes of  $W_\gamma[\check{\nu}]$  over  $\mathbb{Z}/2\mathbb{Z}$  and  $\mathbb{Z}/3\mathbb{Z}$  differ. This is a consequence of the homology of the Klein bottle itself, which depends on the field of coefficients. Over  $\mathbb{Z}/2\mathbb{Z}$ , its first homology group is  $(\mathbb{Z}/2\mathbb{Z})^2$ , while over  $\mathbb{Z}/3\mathbb{Z}$  it is  $\mathbb{Z}/3\mathbb{Z}$ . These homology groups can be read on the right part of the barcodes. In comparison, the barcodes of the usual Čech filtration are the same.

**Example 5.8.** As a last example, we consider two datasets that do not satisfy the hypotheses we studied. Hence the present paper does not provide theoretical guarantees, although our method gives interesting results. The first point cloud, denoted  $X_1$ , is sampled on the unit cube of  $\mathbb{R}^3$ . It is made up of  $6 \times 2000$  points. It can be seen as the immersion of six squares. Note that this immersion does not satisfy the model considered in this paper since the squares are manifolds with boundaries. The second point cloud,  $X_2$ , is sampled on the union of three spheres and a circle. It is made up of  $4 \times 2000$  points. This subset can be seen as the immersion of the disjoint union of three spheres and a circle. Again, this does not fit in our model, since these manifolds have different dimensions. These point clouds are represented in Figure 34.

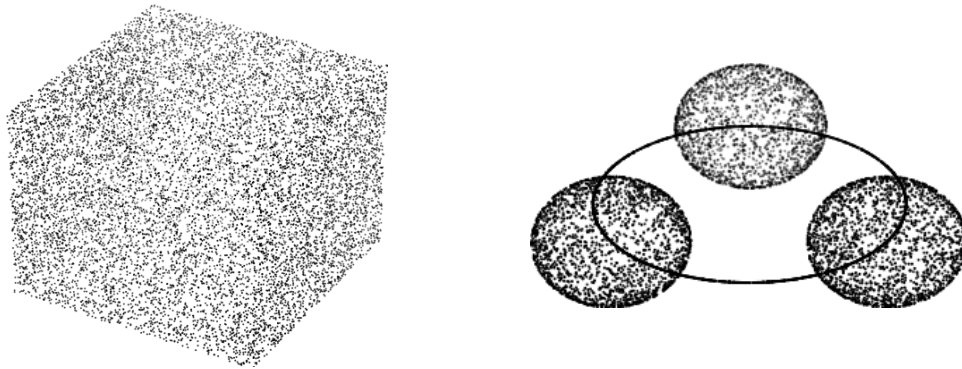


Figure 34: Left:  $X_1$  is a 12'000-sample of the cube. Right:  $X_2$  is a 8'000-sample of the immersion of three spheres and a circle.

We represent on Figure 35 the 0-persistence diagrams of the DTM-filtrations of their lifted measures. We choose the parameters  $r = 0.05$ ,  $\gamma = 2$ ,  $m = 0.01$  for  $X_1$ , and  $r = 0.2$ ,  $\gamma = 2$ ,  $m = 0.01$  for  $X_2$ . Observe that the first barcode contains six long bars, corresponding to the six faces of the cube. Similarly, the second barcode contains four long bars, corresponding to the three spheres and the circle.

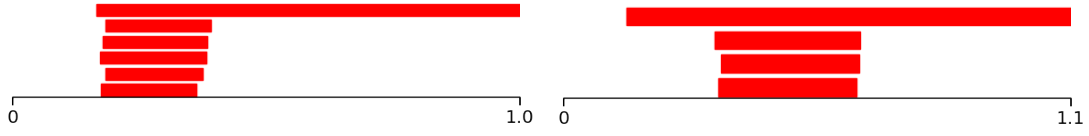


Figure 35: Left: Persistence barcode of the 0-homology of the DTM-filtration of the lifted measure built from  $X_1$ . Right: Same for  $X_2$ .

In Figure 36, we apply the clustering procedure described in Example 5.6. For  $t = 0.35$ ,  $X_1$  is clustered into 83 connected components. We see that there are six main connected components, represented by the faces, and a few outliers. Similarly, we chose  $t = 0.6$  for  $X_2$ , and obtained 20 connected components, four of them representing the four underlying objects.

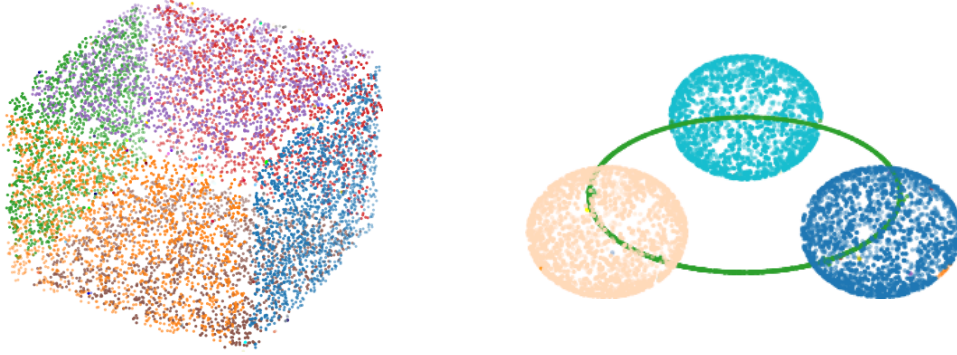


Figure 36: Left: the clustering procedure applied to  $X_1$  at  $t = 0.35$ . Right: Same for  $X_2$  at  $t = 0.6$ .

## 6 Conclusion

In this paper we described a method to estimate the tangent bundle of a manifold  $\mathcal{M}_0$  immersed in a Euclidean space, based on a sample of its image. This estimation is stable in Wasserstein distance. Using the DTM, we are able to estimate the homotopy type of  $\mathcal{M}_0$ . Moreover, via the DTM-filtrations, we can define a filtration of the space  $\mathbb{R}^n \times \mathbf{M}(\mathbb{R}^n)$  whose persistence module contains information about the homology of  $\mathcal{M}_0$ .

The robust estimation of tangent bundles of manifolds opens the way to the estimation of other topological invariants than homology groups—such as characteristic classes—a problem that will be addressed in further works.

Also, as we pointed out in Subsect. 5.2, it would be interesting to understand the geometric quantities associated to the lifted manifold  $\tilde{\mathcal{M}}$  (such as  $\check{\rho}_\gamma$ ,  $\check{L}_{0,\gamma}$ ,  $\check{f}_{\min,\gamma}$  and  $\check{f}_{\max,\gamma}$ ) as a function of those associated with the initial manifold  $\mathcal{M}_0$  ( $\rho$ ,  $L_0$ ,  $f_{\min}$  and  $f_{\max}$ ).

**Acknowledgements.** I would like to thank Frédéric Chazal, Marc Glisse and Théo Lacombe for fruitful discussions and corrections. I also thank the anonymous reviewers for their precious corrections and suggestions.

## A Notations

We adopt the following notations:

- $n, d > 0$  are integers.
- If  $x, y \in \mathbb{R}$ ,  $\min(x, y)$  is the minimum of  $x$  and  $y$ .

- $I$  is the interval  $[0, +\infty)$  or  $[0, T]$  for  $T \geq 0$ .
- $E = \mathbb{R}^n$  is the Euclidean space,  $M(E)$  the vector space of  $n \times n$  matrices,  $\mathcal{G}_d(E)$  the Grassmannian of  $d$ -planes in  $E$ .
- $A$  is a subset of  $E$ ,  $\text{med}(A)$  denotes its medial axis,  $\text{reach}(A)$  its reach. For every  $x \in E$ ,  $\text{dist}(x, A)$  is the distance from  $x$  to  $A$ .
- For  $x, y \in E$ ,  $x \perp y$  denotes the orthogonality of  $x$  and  $y$ .
- If  $x, y \in E$ ,  $x \otimes y = x^t y \in M(E)$  is the outer product, and  $x^{\otimes 2} = x \otimes x$ .
- $\|\cdot\|$  is the Euclidean norm on  $E$  and  $\langle \cdot, \cdot \rangle$  the corresponding inner product,  $\|\cdot\|_F$  the Frobenius norm on  $M(E)$ ,  $\|\cdot\|_\gamma$  the  $\gamma$ -norm on  $E \times M(E)$  (defined in Subsect. 4.1).
- $W_p$  is the  $p$ -Wasserstein distance between measures on  $E$ ,  $W_{p,\gamma}$  is the  $(p, \gamma)$ -Wasserstein distance between measures on  $E \times M(E)$  (defined in Subsect. 4.1).
- $\mathcal{H}^d$  is the  $d$ -dimensional Hausdorff measure on  $E$  or on a subspace  $T \subset E$  (not renormalized).
- If  $\mu$  is a measure of positive finite mass,  $|\mu|$  denotes its mass,  $\bar{\mu} = \frac{1}{|\mu|}\mu$  is the associated probability measure,  $\check{\mu}$  denotes the associated lifted measure (introduced in Subsect. 4.1).
- $1_A$  is the indicator function of a measurable set  $A$ .
- If  $T$  is a subspace of  $E$ ,  $p_T$  denotes the orthogonal projection matrix on  $T$ .
- $\mathcal{B}(x, r)$  and  $\overline{\mathcal{B}}(x, r)$  denote the open and closed balls of  $E$ ,  $\partial \mathcal{B}(x, r)$  the sphere.  $V_d$  and  $S_{d-1}$  denote  $\mathcal{H}^d(\mathcal{B}(0, 1))$  and  $\mathcal{H}^{d-1}(\partial \mathcal{B}(0, 1))$  (note that  $S_{d-1} = dV_d$ ).
- $\mathcal{M}_0$  is a Riemannian manifold, and  $\mathcal{B}_{\mathcal{M}_0}(x, r)$  and  $\overline{\mathcal{B}}_{\mathcal{M}_0}(x, r)$  denote the open and closed geodesics balls. For  $x_0, y_0 \in \mathcal{M}_0$ ,  $d_{\mathcal{M}_0}(x_0, y_0)$  denotes the geodesic distance.
- If  $T$  is a subspace of  $E$ ,  $\mathcal{B}_T(x, r)$  and  $\overline{\mathcal{B}}_T(x, r)$  denote the open and closed balls of  $T$  for the Euclidean distance.
- If  $f$  is a map with values in  $\mathbb{R}$  and  $t \in \mathbb{R}$ ,  $f^t$  denotes the sublevel set  $f^t = f^{-1}((-\infty, t])$ .

## B Table of constants

In the following table, each constant is preceded by the result where it appeared first. If a constant is defined from the others, it is indicated here. The indices are arbitrary and only reflect the order of apparition of each result.



Index	Result	Constant
1.	Corollary 2.9	$a, \quad c_1 = 8\text{diam}(\text{supp}(\mu)) + 5, \quad c'_1 = 2a^{-\frac{1}{d}}$
2.	Hypothesis 2	$\rho$
3.	Hypothesis 3	$L_0, \quad f_{\min}, \quad f_{\max}$
4.	Hypothesis 4	$c_4, \quad r_4$
5.	Lemma 3.4	$c_5: t \mapsto \frac{1}{t} (1 - \sqrt{1-2t})$
6.	Lemma 3.9	$J_{\min} = (\frac{23}{24})^d, \quad J_{\max} = (\frac{5}{4})^d$
7.	Lemma 3.11	$c_7 = 4L_0J_{\max} + \frac{d}{2}\rho f_{\max}$
8.	Proposition 3.13	$c_8 = c_7 + f_{\max}J_{\max}d2^d\rho$
9.	Proposition 3.13, Hypothesis 5	$c_9 = f_{\min}J_{\min}V_d$
10.	Proposition 3.13, Hypothesis 6	$c_{10} = d2^d f_{\max}J_{\max}V_d$
11.	Proposition 3.14, Hypothesis 7	$c_{11} = \frac{f_{\max}J_{\max}}{f_{\min}J_{\min}} \left( \frac{\rho}{\sqrt{4-\sqrt{13}}} \right)^d d2^{2d}\sqrt{3}$
12.	Subsect. 3.3	$\Delta, \quad \Delta_0, \quad \Theta$
13.	Proposition 3.19	$r_{13} = \min \left( \frac{\sin(\Theta)}{8\rho}, \frac{\sin(\Theta)^2}{4}, \frac{\Delta_0 \sin(\Theta)}{4}, \Delta \right)$ $c_{13} = \left( \frac{2}{\sin(\Theta)} \right)^\alpha V_\alpha f_{\max} \mathcal{H}_{\mathcal{M}_0}^{d'}(\mathcal{N}_0)$
14.	Proposition 4.1	$c_{14} = 6\rho + \frac{1}{f_{\min}J_{\min}} (4c_7 + f_{\max}2^d d\rho + c_8),$
15.	Lemma 4.5	$c_{15} = 2 \left( 1 + 4 \frac{5^{d-1}}{3^d} \right) \frac{c_{10}}{c_9}$
16.	Lemma 4.6	$c_{16} = \left( 2 + \frac{2^{\frac{5}{2}} 5^{d-\frac{1}{2}}}{3^d} \right) \frac{c_{11}}{c_9}$
17.	Lemma 4.7	$c_{17} = \frac{2^{d-1}}{c_9} + 2 \frac{12 \cdot 5^{d-1} c_{10} + 1}{3^d c_9} + 2^{d+3} \frac{(\frac{3}{2})^{d-1} c_{10} + 1}{c_9}$
18.	Lemma 4.8	$c_{18} = \frac{2^{d-2}}{c_9} + \frac{4 \cdot 3^{\frac{1}{2}} 5^{d-\frac{1}{2}} c_{11} + 4^{d-\frac{1}{2}}}{3^d c_9} + 2 \cdot 4^d \frac{2c_{11}(\frac{3}{2})^{d-\frac{1}{2}} + 1}{3^d c_9}$
19.	Proposition 4.11	$c_{19} = 4(1 + c_{20}), \quad c'_{19} = 4c_{18}$
20.	Lemma 4.13	$c_{20} = 3 + c_{15} + c_{16} + c_{17}$
21.	Theorem 4.14	$c_{21} = 2 + \frac{1}{2}c'_{19} = 2(1 + c_{18})$
22.	Corollary 4.17	$c_{22} = c_{21}(c_4)^{\frac{1}{p}} + c_{19} + c_{14}$
23.	Subsect. 5.2	$\check{\rho}_\gamma, \quad \check{f}_{\min, \gamma}, \quad \check{c}_{23, \gamma} = \check{f}_{\min, \gamma} J_{\min} V_d$
24.	Corollary 5.5	$\check{c}_{1, \gamma} = 8\text{diam}(\mathcal{M}) + 8\gamma + 5, \quad \check{c}'_{1, \gamma} = 2(\check{c}_{23, \gamma})^{-\frac{1}{d}}$

## C Supplementary material for Sect. 2

*Proof of Lemma 2.4.* The proof is based on the following observations. We can use the triangle inequality, then the Pythagorean Theorem with  $\langle v, y-x \rangle = 0$  and Lemma 2.3 Point 1 to get

$$\begin{aligned}
\|\gamma(t) - x\| &\leq \|(y + tv) - x\| + \|\gamma(t) - (y + tv)\| \\
&\leq \sqrt{\|tv\|^2 + \|y - x\|^2} + \frac{\rho}{2}t^2 \\
&= \sqrt{t^2 + l^2} + \frac{\rho}{2}t^2.
\end{aligned}$$

For any  $r \leq \frac{1}{\rho}$ , consider the equation

$$\sqrt{t^2 + l^2} + \frac{\rho}{2}t^2 = r. \quad (48)$$

By squaring this equality, we get  $(\frac{\rho}{2})^2 t^4 - (1 + \rho r)t^2 + (r^2 - l^2) = 0$ . By considering the polynomial  $T \mapsto (\frac{\rho}{2})^2 T^2 - (1 + \rho r)T + (r^2 - l^2)$ , whose discriminant is  $1 + 2\rho r + (\rho l)^2 > 0$ , we see that the solutions of Equation (48) are

$$T_1 = \frac{\sqrt{2}}{\rho} \sqrt{1 + \rho r - \sqrt{1 + 2\rho r + \rho^2 l^2}} \quad \text{and} \quad T_1' = \frac{\sqrt{2}}{\rho} \sqrt{1 + \rho r + \sqrt{1 + 2\rho r + \rho^2 l^2}}.$$

Following the same ideas, one obtains

$$\|\gamma(t) - x\| \geq \sqrt{t^2 + l^2} - \frac{\rho}{2}t^2.$$

Moreover, the equation

$$\sqrt{t^2 + l^2} - \frac{\rho}{2}t^2 = r \quad (49)$$

admits the following roots:

$$T_2 = \frac{\sqrt{2}}{\rho} \sqrt{1 - \rho r - \sqrt{1 - 2\rho r + \rho^2 l^2}} \quad \text{and} \quad T_2' = \frac{\sqrt{2}}{\rho} \sqrt{1 - \rho r + \sqrt{1 - 2\rho r + \rho^2 l^2}}.$$

We now prove the five points successively.

*Point 1.* Observe that  $\dot{\phi}(t) = 2 \langle \dot{\gamma}(t), \gamma(t) - x \rangle$ , and that

$$\ddot{\phi}(t) = 2 \langle \ddot{\gamma}(t), \gamma(t) \rangle + 2 \langle \dot{\gamma}(t), \gamma(t) - x \rangle.$$

By Cauchy-Schwarz inequality,  $\langle \dot{\gamma}(t), \gamma(t) - x \rangle \geq -\|\dot{\gamma}(t)\| \|\gamma(t) - x\|$ . Note that  $\langle \dot{\gamma}(t), \dot{\gamma}(t) \rangle = 1$  since  $\gamma$  is parametrized by arc-length, and that  $\|\dot{\gamma}(t)\| \leq \rho$  by Equation (1). Hence we get

$$\ddot{\phi}(t) \geq 2(1 - \rho \|\gamma(t) - x\|). \quad (50)$$

Consider Equation (48) with  $r = \frac{1}{\rho}$ . We see that  $\|\gamma(t) - x\| \leq \frac{1}{\rho}$  when  $t$  is lower than

$$T_1 = \frac{\sqrt{2}}{\rho} \sqrt{2 - \sqrt{3 + \rho^2 l^2}}.$$

In this case,  $\ddot{\phi}(t) \geq 0$  according to Equation (50). Since  $\dot{\phi}(0) = 0$ , we deduce that  $\phi$  is increasing on  $[0, T_1]$ .

*Point 2.* As we have seen with Equation (49), we have  $\|\gamma(t) - x\| > r$  when  $t \in (T_2, T_2')$ . In order to give an upper bound on  $T_2$ , we use the inequality  $\sqrt{b} - \sqrt{a} = \frac{1}{\sqrt{a} + \sqrt{b}}(b - a) \leq \frac{1}{\sqrt{b}}(b - a)$ , where  $a < b$ , to get

$$1 - \rho r - \sqrt{1 - 2\rho r + \rho^2 l^2} \leq \frac{1}{1 - \rho r} \rho^2 (r^2 - l^2)$$

and we conclude that  $T_2 \leq \frac{\sqrt{2}}{\sqrt{1 - \rho r}} \sqrt{r^2 - l^2}$ . Since  $r \leq \frac{1}{2\rho}$ , we obtain  $T_2 \leq 2\sqrt{r^2 - l^2}$ .

*Point 3.* When  $l = 0$ , algebraic manipulations show that  $T_2 = \frac{1}{\rho}(1 - \sqrt{1 - 2\rho r})$  and  $T_2' = \frac{1}{\rho}(1 + \sqrt{1 - 2\rho r})$ .

*Point 4.* Consider the map  $\phi: t \mapsto \|\gamma(t) - x\|^2$ . By definition of  $b$ , for all  $t \in (0, b)$ , we have  $\|\gamma(t) - x\| \leq r$ . Hence Equation (50) gives  $\ddot{\phi}(t) \geq 2(1 - \rho r)$ . It follows that  $\dot{\phi}(t) \geq 2(1 - \rho r)t$ , and that

$$\begin{aligned} \phi(b) - \phi(a) &= \int_a^b \dot{\phi}(t) dt \geq \int_a^b 2(1 - \rho r)t dt \\ &= (1 - \rho r)(b^2 - a^2). \end{aligned}$$

Note that  $r^2 = \phi(b)$ . Besides,  $s^2 = \phi(a)$  or  $s^2 < \phi(a)$ , depending on whether  $s \geq l$  or  $s < l$ . In both cases, we have  $r^2 - s^2 \geq \phi(b) - \phi(a)$ , and we deduce that

$$r^2 - s^2 \geq (1 - \rho r)(b^2 - a^2).$$

Writing  $r^2 - s^2 = (r + s)(r - s)$  and  $b^2 - a^2 = (b + a)(b - a)$  leads to

$$b - a \leq \frac{r + s}{b + a} \frac{1}{1 - \rho r} (r - s). \quad (51)$$

Now, let us give a lower bound on  $b$ . According to Equation (48),  $b$  is lower bounded by  $T_1 = \frac{\sqrt{2}}{\rho} \sqrt{1 + \rho r - \sqrt{1 + 2\rho r + \rho^2 l^2}}$ . Using the inequality  $\sqrt{b} - \sqrt{a} = \frac{1}{\sqrt{b} + \sqrt{a}}(b - a) \geq \frac{1}{2\sqrt{b}}(b - a)$ , where  $a < b$ , we get

$$1 + \rho r - \sqrt{1 + 2\rho r + \rho^2 l^2} \geq \frac{1}{2(1 + \rho r)} \rho^2 (r^2 - l^2),$$

and we conclude that  $b \geq (1 + \rho r)^{-\frac{1}{2}} \sqrt{r^2 - s^2}$ . Injecting  $b + a \geq b \geq (1 + \rho r)^{-\frac{1}{2}} \sqrt{r^2 - s^2}$  in Equation (51) yields

$$b(v) - a(v) \leq \frac{(1 + \rho r)^{\frac{1}{2}}}{1 - \rho r} \sqrt{r^2 - s^2}.$$

Under the hypothesis  $r \leq \frac{1}{2\rho}$ , we get  $b - a \leq \sqrt{6} \sqrt{r^2 - s^2}$ .

*Point 5.* When  $l = 0$ , we have  $b(v) + a(v) \geq r + s$ . Hence Equation (51) yields  $b(v) - a(v) \leq \frac{1}{1 - \rho r} (r - s)$ . Using  $r \leq \frac{1}{2\rho}$ , we obtain  $b(v) - a(v) \leq 2(r - s)$ .  $\square$

*of Corollary 2.9.* We shall first study an intermediate quantity. Let  $\mu$  be a probability measure on  $E = \mathbb{R}^n$ ,  $m \in (0, 1)$ , and  $d_{\mu, m}$  the corresponding DTM. Consider the quantity  $c(\mu, m)$  is defined as

$$c(\mu, m) = \sup_{x \in \text{supp}(\mu)} d_{\mu, m}(x).$$

Suppose that  $\mu$  satisfies the following for  $r < (\frac{m}{a})^{\frac{1}{d}}$ :  $\forall x \in \text{supp}(\mu), \mu(\mathcal{B}(x, r)) \geq ar^d$ . Let us show that  $c(\mu, m) \leq Cm^{\frac{1}{d}}$  with  $C = a^{-\frac{1}{d}}$ . By definition,

$$\delta_{\mu, t}(x) = \inf \{r \geq 0 \mid \mu(\overline{\mathcal{B}}(x, r)) > t\} \quad \text{and} \quad d_{\mu, m}^2(x) = \frac{1}{m} \int_0^m \delta_{\mu, t}^2(x) dt.$$

Using the assumption  $\mu(\mathcal{B}(x, r)) \geq ar^d$  for all  $x \in \text{supp}(\mu)$ , we get  $\delta_{\mu, t}(x) \leq (\frac{t}{a})^{\frac{1}{d}}$ , and a simple computation yields

$$d_{\mu, m}^2(x) \leq \frac{d}{d+2} \left(\frac{t}{a}\right)^{\frac{2}{d}} \leq \left(\frac{t}{a}\right)^{\frac{2}{d}},$$

which yields the result.

We can now prove the corollary. Let  $\pi$  be an optimal transport plan for  $w = W_2(\mu, \nu)$ . Denote  $\alpha = w^{\frac{1}{2}}$  and  $D = \text{diam}(\text{supp}(\mu))$ . Define  $\pi'$  to be  $\pi$  restricted to the set  $\{x, y \in E \mid \|x - y\| < \alpha\}$ . We denote its marginals  $\mu'$  and  $\nu'$ . By Markov inequality,  $1 - |\pi'| \leq \frac{w^2}{\alpha^2} = w$ , where we recall that  $|\pi'|$  denotes the total mass of  $\pi'$ . Consider the probability measures  $\overline{\mu}' = \frac{1}{|\mu'|} \mu'$  and  $\overline{\nu}' = \frac{1}{|\nu'|} \nu'$ . Let us show that we have

$$W_2(\mu, \overline{\mu}') = 2D\alpha, \quad W_2(\overline{\mu}', \overline{\nu}') \leq \alpha \quad \text{and} \quad W_2(\nu, \overline{\nu}') \leq 2(1 + D)\alpha. \quad (52)$$

The first inequality is an application of Lemma 4.4:

$$W_2(\mu, \overline{\mu}') \leq 2(1 - |\mu'|)^{\frac{1}{2}} D = 2(1 - |\pi'|)^{\frac{1}{2}} D \leq 2w^{\frac{1}{2}} D.$$

To obtain the second inequality, we write

$$\begin{aligned} W_2^2(\bar{\mu}', \bar{\nu}') &= \int \|x - y\|^2 d\bar{\pi}'(x, y) = \int \|x - y\| \frac{d\bar{\pi}'(x, y)}{|\bar{\pi}'|} \\ &\leq \frac{1}{|\bar{\pi}'|} \int \|x - y\| d\bar{\pi}(x, y). \end{aligned}$$

Hence Jensen inequality leads to  $W_2(\bar{\mu}', \bar{\nu}') \leq \frac{w}{|\bar{\pi}'|^{\frac{1}{2}}}$ . Since  $1 - |\bar{\pi}'| \leq w$ , we have  $\frac{w}{|\bar{\pi}'|^{\frac{1}{2}}} \leq \frac{w}{1-w}$ , and the assumption  $w \leq \frac{1}{4}$  yields  $\frac{w}{1-w} \leq \alpha$ . This proves the second point. Finally, we obtain the third inequality by applying the triangle inequality:

$$W_2(\nu, \bar{\nu}') \leq W_2(\nu, \mu) + W_2(\mu, \bar{\mu}') + W_2(\bar{\mu}', \bar{\nu}').$$

Next, let us deduce that

$$\begin{aligned} c(\bar{\mu}', m) &\leq c(\mu) + m^{-\frac{1}{2}} 2D\alpha \\ \text{and} \quad c(\bar{\nu}', m) &\leq c(\mu, m) + \left(m^{-\frac{1}{2}} + m^{-\frac{1}{2}} 2D + 1\right) \alpha. \end{aligned} \tag{53}$$

The first inequality follows from the stability of the DTM (see Equation (3)):

$$c(\bar{\mu}', m) = \sup_{x \in \text{supp}(\bar{\mu}')} d_{\bar{\mu}'}(x) \leq \sup_{x \in \text{supp}(\bar{\mu}')} d_{\mu}(x) + m^{-\frac{1}{2}} W_2(\bar{\mu}', \mu),$$

and we conclude with  $W_2(\mu, \bar{\mu}') = 2D\alpha$ . In order to prove the second inequality, we also use Equation (3):

$$c(\bar{\nu}', m) = \sup_{x \in \text{supp}(\bar{\nu}')} d_{\bar{\nu}'}(x) \leq \sup_{x \in \text{supp}(\bar{\nu}')} d_{\bar{\mu}'}(x) + m^{-\frac{1}{2}} W_2(\bar{\mu}', \bar{\nu}').$$

Since  $\pi'$  has support included in  $\{x, y \in E \mid \|x - y\| < \alpha\}$ , we use the fact that the DTM is 1-Lipschitz to obtain

$$\sup_{x \in \text{supp}(\bar{\nu}')} d_{\bar{\mu}'}(x) \leq \sup_{x \in \text{supp}(\bar{\mu}')} d_{\bar{\mu}'}(x) + \alpha = c(\mu', m) + \alpha$$

and we deduce

$$\begin{aligned} c(\bar{\nu}', m) &\leq c(\mu', m) + \alpha + m^{-\frac{1}{2}} W_2(\bar{\mu}', \bar{\nu}') \\ &\leq c(\mu, m) + (m^{-\frac{1}{2}} + m^{-\frac{1}{2}} 2D + 1) \alpha. \end{aligned}$$

We can now conclude with Theorem 2.8. In our context, it reads

$$\begin{aligned} d_i(W[\mu], W[\nu]) &\leq m^{-\frac{1}{2}} W_2(\mu, \bar{\mu}') + m^{-\frac{1}{2}} W_2(\bar{\mu}', \bar{\nu}') + m^{-1} W_2(\nu, \bar{\nu}') + c(\bar{\mu}', m) + c(\bar{\nu}', m) \\ &\leq (m^{-\frac{1}{2}} (4D + 1) + 4(D + 1)) \alpha + 2c(\mu, m), \end{aligned}$$

where we used Equations (52) and (53) on the last line. Since  $m \leq 1$ , we can simplify this expression into

$$d_i(W[\mu], W[\nu]) \leq m^{-\frac{1}{2}} (8D + 5) \alpha + 2c(\mu, m).$$

We conclude the proof by using the inequality  $c(\mu, m) \leq a^{-\frac{1}{d}} m^{\frac{1}{d}}$  shown at the beginning of the proof.  $\square$

## References

- [1] Allen Hatcher. *Algebraic Topology*. Cambridge University Press, 2002.
- [2] Partha Niyogi, Stephen Smale, and Shmuel Weinberger. Finding the homology of submanifolds with high confidence from random samples. *Discrete & Computational Geometry*, 39(1-3):419–441, 2008.
- [3] Frédéric Chazal and André Lieutier. Smooth manifold reconstruction from noisy and non-uniform approximation with guarantees. *Computational Geometry*, 40(2):156–170, 2008.
- [4] Jisu Kim, Jaehyeok Shin, Frédéric Chazal, Alessandro Rinaldo, and Larry Wasserman. Homotopy reconstruction via the cech complex and the vietoris-rips complex. *arXiv preprint arXiv:1903.06955*, 2019.
- [5] Frédéric Chazal, David Cohen-Steiner, and André Lieutier. A sampling theory for compact sets in Euclidean space. *Discrete & Computational Geometry*, 41(3):461–479, 2009.
- [6] Dominique Attali, André Lieutier, and David Salinas. Vietoris–rips complexes also provide topologically correct reconstructions of sampled shapes. *Computational Geometry*, 46(4):448–465, 2013.
- [7] Sara Kalisnik and Davorin Lesnik. Finding the homology of manifolds using ellipsoids. *arXiv preprint arXiv:2006.09194*, 2020.
- [8] Herbert Edelsbrunner. The union of balls and its dual shape. In *Proceedings of the ninth annual symposium on Computational geometry*, pages 218–231, 1993.
- [9] Herbert Edelsbrunner and Ernst P Mücke. Three-dimensional alpha shapes. *ACM Transactions on Graphics (TOG)*, 13(1):43–72, 1994.
- [10] Vin De Silva and Gunnar E Carlsson. Topological estimation using witness complexes. In *PBG*, pages 157–166, 2004.
- [11] Dominique Attali, Herbert Edelsbrunner, and Yuriy Mileyko. Weak witnesses for delaunay triangulations of submanifolds. In *Proceedings of the 2007 ACM symposium on Solid and physical modeling*, pages 143–150, 2007.
- [12] Jean-Daniel Boissonnat and Arijit Ghosh. Manifold reconstruction using tangential Delaunay complexes. *Discrete & Computational Geometry*, 51(1):221–267, 2014.
- [13] Frédéric Chazal and André Lieutier. Stability and computation of topological invariants of solids in  $\mathbb{R}^n$ . *Discrete & Computational Geometry*, 37(4):601–617, 2007.
- [14] Frédéric Chazal and Steve Yann Oudot. Towards persistence-based reconstruction in Euclidean spaces. In *Proceedings of the twenty-fourth annual symposium on Computational geometry*, SCG ’08, pages 232–241, New York, NY, USA, 2008. ACM.
- [15] Brittany Terese Fasy, Rafal Komendarczyk, Sushovan Majhi, and Carola Wenk. On the reconstruction of geodesic subspaces of  $\mathbb{R}^n$ . *arXiv preprint arXiv:1810.10144*, 2018.
- [16] Herbert Edelsbrunner and John Harer. *Computational topology: an introduction*. American Mathematical Soc., 2010.
- [17] Jean-Daniel Boissonnat, Frédéric Chazal, and Mariette Yvinec. *Geometric and topological inference*, volume 57. Cambridge University Press, 2018.
- [18] Frédéric Chazal, Vin de Silva, Marc Glisse, and Steve Oudot. *The Structure and Stability of Persistence Modules*. SpringerBriefs in Mathematics, 2016.

- [19] Peter Bubenik, Gunnar Carlsson, Peter T Kim, and Zhi-Ming Luo. Statistical topology via morse theory persistence and nonparametric estimation. *Algebraic methods in statistics and probability II*, 516:75–92, 2010.
- [20] Brittany Terese Fasy, Fabrizio Lecci, Alessandro Rinaldo, Larry Wasserman, Sivaraman Balakrishnan, and Aarti Singh. Confidence sets for persistence diagrams. *The Annals of Statistics*, 42(6):2301–2339, 2014.
- [21] Katharine Turner, Yuriy Mileyko, Sayan Mukherjee, and John Harer. Fréchet means for distributions of persistence diagrams. *Discrete & Computational Geometry*, 52(1):44–70, 2014.
- [22] F. Chazal, D. Cohen-Steiner, and Q. Mérigot. Geometric inference for probability measures. *Journal on Found. of Comp. Mathematics*, 11(6):733–751, 2011.
- [23] J. Phillips, B. Wang, and Y Zheng. Geometric inference on kernel density estimates. In *Proc. 31st Annu. Sympos. Comput. Geom (SoCG 2015)*, pages 857–871, 2015.
- [24] Leonidas Guibas, Dmitriy Morozov, and Quentin Mérigot. Witnessed k-distance. *Discrete & Computational Geometry*, 49(1):22–45, 2013.
- [25] Mickaël Buchet, Frédéric Chazal, Steve Y Oudot, and Donald R Sheehy. Efficient and robust persistent homology for measures. *Computational Geometry*, 58:70–96, 2016.
- [26] Hirokazu Anai, Frédéric Chazal, Marc Glisse, Yuichi Ike, Hiroya Inakoshi, Raphaël Tinarrage, and Yuhei Umeda. DTM-based filtrations. In *Topological Data Analysis*, pages 33–66. Springer, 2020.
- [27] Yong Wang, Yuan Jiang, Yi Wu, and Zhi-Hua Zhou. Spectral clustering on multiple manifolds. *IEEE Transactions on Neural Networks*, 22(7):1149–1161, 2011.
- [28] Dian Gong, Xuemei Zhao, and Gérard Medioni. Robust multiple manifolds structure learning. *arXiv preprint arXiv:1206.4624*, 2012.
- [29] Ery Arias-Castro, Gilad Lerman, and Teng Zhang. Spectral clustering based on local PCA. *The Journal of Machine Learning Research*, 18(1):253–309, 2017.
- [30] Siu-Wing Cheng and Man-Kwun Chiu. Tangent estimation from point samples. *Discrete & Computational Geometry*, 56(3):505–557, 2016.
- [31] Eddie Aamari, Jisu Kim, Frédéric Chazal, Bertrand Michel, Alessandro Rinaldo, and Larry Wasserman. Estimating the Reach of a Manifold. *Electronic journal of statistics*, 2019.
- [32] Eddie Aamari and Clément Levrard. Nonasymptotic rates for manifold, tangent space and curvature estimation. *The Annals of Statistics*, 47(1):177–204, 2019.
- [33] Amit Singer and H-T Wu. Vector diffusion maps and the connection laplacian. *Communications on pure and applied mathematics*, 65(8):1067–1144, 2012.
- [34] Jisu Kim, Alessandro Rinaldo, and Larry Wasserman. Minimax rates for estimating the dimension of a manifold. *arXiv preprint arXiv:1605.01011*, 2016.
- [35] Vladimir I Koltchinskii. Empirical geometry of multivariate data: a deconvolution approach. *Annals of statistics*, pages 591–629, 2000.
- [36] Anna V Little, Jason Lee, Yoon-Mo Jung, and Mauro Maggioni. Estimation of intrinsic dimensionality of samples from noisy low-dimensional manifolds in high dimensions with multiscale SVD. In *2009 IEEE/SP 15th Workshop on Statistical Signal Processing*, pages 85–88. IEEE, 2009.
- [37] Philippos Mordohai and Gérard Medioni. Dimensionality estimation, manifold learning and function approximation using tensor voting. *Journal of Machine Learning Research*, 11(1), 2010.

- [38] Diego H Díaz Martínez, Facundo Mémoli, and Washington Mio. The shape of data and probability measures. *Applied and Computational Harmonic Analysis*, 2018.
- [39] Facundo Memoli, Zane Smith, and Zhengchao Wan. The Wasserstein transform. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 4496–4504, Long Beach, California, USA, 09–15 Jun 2019. PMLR.
- [40] Blanche Buet, Gian Paolo Leonardi, and Simon Masnou. A varifold approach to surface approximation. *Archive for Rational Mechanics and Analysis*, 226(2):639–694, 2017.
- [41] Blanche Buet, Gian Paolo Leonardi, and Simon Masnou. Weak and approximate curvatures of a measure: a varifold perspective. *arXiv preprint arXiv:1904.05930*, 2019.
- [42] Nicolas Charon and Alain Trounev. The varifold representation of nonoriented shapes for diffeomorphic registration. *SIAM Journal on Imaging Sciences*, 6(4):2547–2580, 2013.
- [43] Herbert Federer. Curvature measures. *Transactions of the American Mathematical Society*, 93(3):418–491, 1959.
- [44] M.P. do Carmo. *Riemannian Geometry*. Mathematics (Boston, Mass.). Birkhäuser, 1992.
- [45] Frank Morgan. *Geometric measure theory: a beginner’s guide*. Academic press, 2016.
- [46] Alexander Lytchak. Almost convex subsets. *Geometriae Dedicata*, 115(1):201–218, 2005.
- [47] Jean-Daniel Boissonnat, André Lieutier, and Mathijs Wintraecken. The reach, metric distortion, geodesic convexity and the variation of tangent spaces. *Journal of Applied and Computational Topology*, 3(1-2):29–58, 2019.
- [48] Stephanie B Alexander and Richard L Bishop. Gauss equation and injectivity radii for subspaces in spaces of curvature bounded above. *Geometriae Dedicata*, 117(1):65–84, 2006.
- [49] Eddie Aamari. *Vitesses de convergence en inférence géométrique*. PhD thesis, Université Paris-Saclay, 2018.
- [50] Magnus Botnan and William Crawley-Boevey. Decomposition of persistence modules. *Proceedings of the American Mathematical Society*, 148(11):4581–4596, 2020.
- [51] Ralph J Herbert. *Multiple points of immersed manifolds*, volume 250. American Mathematical Soc., 1981.
- [52] Alfred Gray. *Tubes*, volume 221. Birkhäuser, 2012.
- [53] Raphaël Tinarrage. Computing persistent Stiefel–Whitney classes of line bundles. *Journal of Applied and Computational Topology*, pages 1–61, 2021.