

# WITNESSED K-DISTANCE

LEONIDAS GUIBAS, DMITRIY MOROZOV, AND QUENTIN MÉRIGOT

ABSTRACT. Distance function to a compact set plays a central role in several areas of computational geometry. Methods that rely on it are robust to the perturbations of the data by the Hausdorff noise, but fail in the presence of outliers. The recently introduced *distance to a measure* offers a solution by extending the distance function framework to reasoning about the geometry of probability measures, while maintaining theoretical guarantees about the quality of the inferred information. A combinatorial explosion hinders working with distance to a measure as an ordinary power distance function. In this paper, we analyze an approximation scheme that keeps the representation linear in the size of the input, while maintaining the guarantees on the inference quality close to those for the exact but costly representation.

## 1. INTRODUCTION

The problem of recovering the geometry and topology of compact sets from finite point samples has seen several important developments in the previous decade. Homeomorphic surface reconstruction algorithms have been proposed to deal with surfaces in  $\mathbb{R}^3$  sampled without noise [1] and with moderate Hausdorff (local) noise [13]. In the case of submanifolds of a higher dimensional Euclidean space [21], or even for more general compact subsets [5], it is also possible, at least in principle, to compute the homotopy type from a Hausdorff sampling. If one is only interested in the homology of the underlying space, the theory of persistent homology [16] applied to Vietoris-Rips complexes provides an algorithmically tractable way to estimate the Betti numbers from a finite Hausdorff sampling [7].

All of these constructions share a common feature: they estimate the geometry of the underlying space by a union of balls of some radius  $r$  centered at the data points  $P$ . A different interpretation of this union is the  $r$ -sublevel set of the *distance function* to  $P$ ,  $d_P : x \mapsto \min_{p \in P} \|x - p\|$ . Distance functions capture the geometry of their defining sets, and they are stable to Hausdorff perturbations of those sets, making them well-suited for reconstruction results. However, they are also extremely sensitive to the presence of outliers (i.e. data points that lie far from the underlying set); all reconstruction techniques that rely on them fail even in presence of a single outlier.

To counter this problem, Chazal, Cohen-Steiner, and Mériqot [6] developed a notion of *distance function to a probability measure* that retains the properties of the (usual) distance important for geometric inference. Instead of assuming an underlying compact set that is sampled by the points, they assume an underlying probability measure  $\mu$  from which the point sample  $P$  is drawn. The distance function  $d_{\mu, m_0}$  to the measure  $\mu$  depends on a mass parameter  $m_0 \in (0, 1)$ . This parameter acts as a smoothing term: a smaller  $m_0$  captures the geometry of the support better, while a larger  $m_0$  leads to better stability at the price of precision.

Crucially, the function  $d_{\mu, m_0}$  is stable to the perturbations of the measure  $\mu$  under the Wasserstein distance, defined in Section 2.2. For instance, the Wasserstein distance between the underlying measure  $\mu$  and the uniform probability measure on the point set  $P$  can be small even if  $P$  contains some outliers. When this happens, the stability result ensures that distance function  $d_{\mathbf{1}_P, m_0}$  to the uniform probability measure  $\mathbf{1}_P$  on  $P$  retains the geometric information contained in the underlying measure  $\mu$  and its support.

**Computing with distance functions to measures.** In this article we address the computational issues related to this new notion. If  $P$  is a subset of  $\mathbb{R}^d$  containing  $N$  points, and  $m_0 = k/N$ , we will denote the distance function to the uniform measure on  $P$  by  $d_{P,k}$ . As observed in [6], the value of  $d_{P,k}$  at a given point  $x$  is easy to compute: it is the square root of the average squared distance from the point  $x$  to its  $k$  nearest neighbors in  $P$ . However, most inference methods require a way to represent the function — more precisely, its sublevel sets — globally. It turns out that the distance function  $d_{P,k}$  can be rewritten as a minimum

$$(1) \quad d_{P,k}^2(x) = \min_{\bar{p}} \|x - \bar{p}\|^2 - w_{\bar{p}},$$

where  $\bar{p}$  ranges over the set of barycenters of  $k$  points in  $P$  (see Section 3). Computational geometry provides a rich toolbox to represent sublevel sets of such functions, for example, via weighted  $\alpha$ -complexes [15].

The difficulty in applying these methods is that to get an equality in Equation (1) the minimum number of barycenters to store is the same as the number of sites in the order- $k$  Voronoi diagram of  $P$ , making this representation unusable even for modest input sizes [9]. Our solution is to construct an approximation of the distance function  $d_{P,k}$ , defined by the same equation as (1), but with  $\bar{p}$  ranging over a smaller subset of barycenters. In this article, we study the quality of approximation given by a *linear-sized* subset — the *witnessed barycenters*, defined as the barycenters of any  $k$  points in  $P$  whose order- $k$  Voronoi cell contains at least one of the sample points. The algorithmic simplicity of the scheme is appealing: we only have to find the  $k - 1$  nearest neighbors for each input point. We denote by  $d_{P,k}^w$  and call *witnessed  $k$ -distance* the function defined by Equation (1), where  $\bar{p}$  ranges over the witnessed barycenters.

**Contributions.** Our goal is to give conditions on the point cloud  $P$  under which the witnessed  $k$ -distance  $d_{P,k}^w$  provides a good uniform approximation of the distance to measure  $d_{P,k}$ . We first give a general multiplicative bound on the error produced by this approximation. However, most of our paper (Sections 4 and 5) analyzes the uniform approximation error, when  $P$  is a set of independent samples from a measure concentrated near a lower-dimensional subset of the Euclidean space. The following is a prototypical example for this setting, although our analysis allows for a wider range of problems. Note that some of the common settings in the literature fit either directly into this example, or into its logic — for instance, the mixture of Gaussians [12] and off-manifold Gaussian noise in normal directions [20].

- (H1) We assume that the “ground truth” is an unknown probability measure  $\mu$  whose *dimension is bounded* by a small constant  $\ell$ . Practically, this means that  $\mu$  is concentrated on a compact set  $K \subseteq \mathbb{R}^d$  whose dimension is at most  $\ell$ , and whose mass distribution “remembers” all of  $K$  (see Definition 3). A

prototypical  $\mu$  is the uniform measure on a smooth compact  $\ell$ -dimensional submanifold  $K$ , or on a finite union of such submanifolds.

This hypothesis ensures that the distance to the measure  $\mu$  is close to the distance to the support  $K$  of  $\mu$ , and lets us recover information about  $K$ . Our first result (Witnessed Bound Theorem 2) states that if the uniform measure to a point cloud  $P$  is a good Wasserstein-approximation of  $\mu$ , then the witnessed  $k$ -distance to  $P$  provides a good approximation of the distance to the underlying compact set  $K$ . Our bound is only a constant times worse than the bound for the exact  $k$ -distance.

- (H2) The second assumption is that we are not sampling directly from  $\mu$ , but through a noisy channel — another measure  $\nu$ . For instance,  $\nu$  could result from the convolution of  $\mu$  with a Gaussian distribution  $\mathcal{N}(0, d^{-1}\sigma^2\mathbf{I})$  whose variance is  $\sigma^2$ . The precise condition on  $\nu$  is given in Definition 5; it introduces the notion of  $\sigma$ -perturbation, related to the Wasserstein distance between the measures  $\mu$  and  $\nu$ . This generalization allows, in particular, to consider noise models that are not translation-invariant.
- (H3) Finally, we suppose that our input data set  $P \subseteq \mathbb{R}^d$  consists of  $N$  points drawn independently from the noisy measure  $\nu$ . We denote with  $\mathbf{1}_P$  the uniform measure on  $P$ .

These two hypotheses let us control the Wasserstein distance between  $\mu$  and  $\mathbf{1}_P$  with high probability. We assume that the point cloud  $P$  is gathered following the three hypothesis above. Our second result states that the witnessed  $k$ -distance to  $P$  provides a good approximation of the distance to the compact set  $K$  with high probability, as soon as the amount of noise  $\sigma$  is low enough and the number of points  $N$  is large enough.

APPROXIMATION THEOREM (THEOREM 4). Let  $P$  be a set of  $N$  points drawn according to the three hypotheses (H1)-(H3), let  $k \in \{1, \dots, N\}$  and  $m_0 = k/N$ . Then, the error bound

$$\|d_{P,k}^w - d_K\|_\infty \leq 18m_0^{-1/2}\sigma + 12m_0^{1/\ell}\alpha_\mu^{-1/\ell}$$

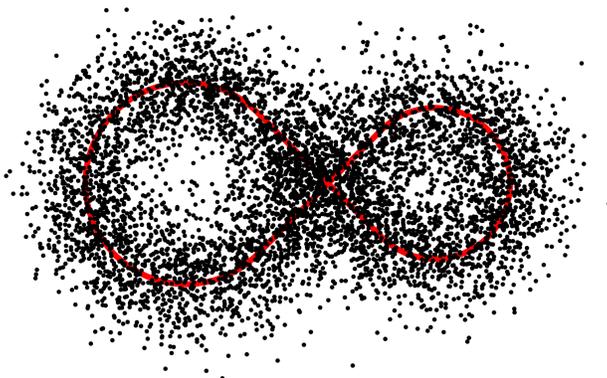
holds with high probability.

**Outline.** The relevant background appears in Section 2. We present our approximation scheme together with a general bound of its quality in Section 3. We analyze its approximation quality for measures concentrated on low-dimensional subsets of the Euclidean space in Section 4. The convergence of the uniform measure on a point cloud sampled from a measure of low complexity appears in Section 5 and leads to our main result. We illustrate the utility of the bound with an example and a topological inference statement in our final Section 6.

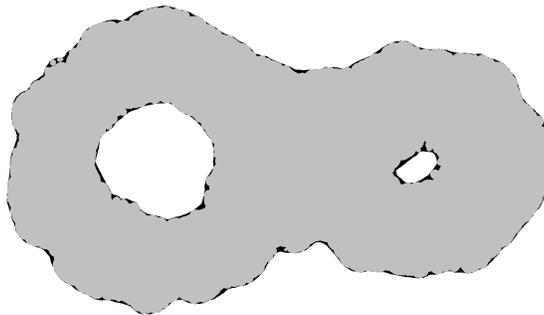
## 2. BACKGROUND

We begin by reviewing the relevant background.

**2.1. Measure.** Let us briefly recap the few concepts of measure theory that we use. A *non-negative measure*  $\mu$  on the space  $\mathbb{R}^d$  is a map from (Borel) subsets of  $\mathbb{R}^d$  to non-negative numbers, which is *additive* in the sense that  $\mu(\cup_{i \in \mathcal{N}} B_i) = \sum_i \mu(B_i)$  whenever  $(B_i)$  is a countable family of disjoint (Borel) subsets. The *total mass* of a measure  $\mu$  is  $\text{mass}(\mu) := \mu(\mathbb{R}^d)$ . A measure  $\mu$  with unit total mass is called a *probability measure*. The *support* of a measure  $\mu$ , denoted by  $\text{spt}(\mu)$ ,



(a) Data



(b) Sublevel sets

FIGURE 1. (a) 6000 points sampled from a sideways figure 8 (in red), with circle radii  $R_1 = \sqrt{2}$  and  $R_2 = \sqrt{9/8}$ . The points are sampled from the uniform measure on the figure-8, convolved with the Gaussian distribution  $\mathcal{N}(0, \sigma^2)$ , where  $\sigma = .45$ . (b)  $r$ -sublevel sets of the witnessed (in gray) and exact (additional points in black)  $k$ -distances with mass parameter  $m_0 = 50/6000$ , and  $r = .239$ .

is the smallest closed set whose complement has zero measure. The *expectation* or *mean* of  $\mu$  is the point  $\mathbb{E}(\mu) = \int_{\mathbb{R}^d} x d\mu(x)$ ; the variance of  $\mu$  is the number  $\sigma_\mu^2 = \int_{\mathbb{R}^d} \|x - \mathbb{E}(\mu)\|^2 d\mu(x)$ .

Although the results we present are often more general, the typical probability measures we have in mind are of two kinds: (i) the probability measure defined by rescaling the volume form of a lower-dimensional submanifold of the ambient space and (ii) discrete probability measures that are obtained through noisy sampling of probability measures of the previous kind. For any finite set  $P$  with  $N$  points, denote by  $\mathbf{1}_P$  the uniform measure supported on  $P$ , i.e. the sum of Dirac masses centered at  $p \in P$  with weight  $1/N$ .

**2.2. Wasserstein distance.** A natural way to quantify the distance between two measures is the *Wasserstein distance*. This distance measures the  $L^2$ -cost of transporting the mass of the first measure onto the second one. A general study of this notion and its relation to the problem of optimal transport appear in [22]. We first give the general definition and then explain its interpretation when one of the two measures has finite support.

A *transport plan* between two measures  $\mu$  and  $\nu$  with the same total mass is a measure  $\pi$  on the product space  $\mathbb{R}^d \times \mathbb{R}^d$  such that for every subsets  $A, B$  of  $\mathbb{R}^d$ ,  $\pi(A \times \mathbb{R}^d) = \mu(A)$  and  $\pi(\mathbb{R}^d \times B) = \nu(B)$ . Intuitively,  $\pi(A \times B)$  represents the amount of mass of  $\mu$  contained in  $A$  that will be transported to  $B$  by  $\pi$ . The *cost* of this transport plan is given by

$$c(\pi) := \left( \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^2 d\pi(x, y) \right)^{1/2}$$

Finally, the *Wasserstein distance* between  $\mu$  and  $\nu$  is the minimum cost of a transport plan between these measures.

Consider the special case where the measure  $\nu$  is supported on a finite set  $P$ . This means that  $\nu$  can be written as  $\sum_{p \in P} \alpha_p \delta_p$ , where  $\delta_p$  is the unit Dirac mass at  $p$ . Moreover,  $\sum_p \alpha_p$  must equal the total mass of  $\mu$ . A transport plan  $\pi$  between  $\mu$  and  $\nu$  corresponds to a decomposition of  $\mu$  into a sum of positive measures  $\sum_{p \in P} \mu_p$  such that  $\text{mass}(\mu_p) = \alpha_p$ . The squared cost of the plan defined by this decomposition is then

$$c(\pi) = \left( \sum_{p \in P} \left[ \int_{\mathbb{R}^d} \|x - p\|^2 d\mu_p(x) \right] \right)^{1/2}.$$

**Wasserstein noise.** Two properties of the Wasserstein distances are particularly useful to us. Together, they show that the Wasserstein noise and sampling model generalize the commonly used empirical sampling with Gaussian noise model:

- Consider a probability measure  $\mu$  and  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  the density of a probability distribution centered at the origin, and denote by  $\nu$  the result of the convolution of  $\mu$  by  $f$ . Then, the Wasserstein distance between  $\mu$  and  $\nu$  is at most  $\sigma$ , where  $\sigma^2 := \int_{\mathbb{R}^d} \|x\|^2 f(x) dx$  is the variance of the probability distribution defined by  $f$ .
- Let  $P$  denote a set of  $N$  points drawn independently from a given measure  $\nu$ . Then, the Wasserstein distance  $W_2(\nu, \mathbf{1}_P)$  between  $\nu$  and the uniform probability measure on  $P$  converges to zero as  $N$  grows to infinity with high probability. Examples of such asymptotic convergence results are common in statistics, e.g. [4] and references therein. In Theorem 3 below, we give a quantitative non-asymptotic result assuming that  $\nu$  is low-dimensional **(H1)**.

Using the notation introduced in the two items above, one can show that with high probability  $\limsup_N W_2(\mathbf{1}_P, \mu) \leq \sigma$ . A more quantitative version of this statement appears in Corollary 1 below.

**2.3. Distance-to-measure and  $k$ -distance.** In [6], the authors introduce a distance to a probability measure as a way to infer the geometry and topology of this

measure in the same way the geometry and topology of a set is inferred from its distance function. Given a probability measure  $\mu$  and a *mass parameter*  $m_0 \in (0, 1)$ , they define a distance function  $d_{\mu, m_0}$  which captures the properties of the usual distance function to a compact set that are used for geometric inference.

DEFINITION 1. For any point  $x$  in  $\mathbb{R}^d$ , let  $\delta_{\mu, m}(x)$  be the radius of the smallest ball centered at  $x$  that contains a mass at least  $m$  of the measure  $\mu$ . The *distance to the measure*  $\mu$  with parameter  $m_0$  is defined by  $d_{\mu, m_0}(x) = m_0^{-1/2} \left( \int_{m=0}^{m_0} \delta_{\mu, m}(x)^2 dm \right)^{1/2}$ .

Given a point cloud  $P$  containing  $N$  points, the measure of interest is the uniform measure  $\mathbf{1}_P$  on  $P$ . When  $m_0$  is a fraction  $k/N$  of the number of points (where  $k$  is an integer), we call *k-distance* and denote by  $d_{P, k}$  the distance to the measure  $d_{\mathbf{1}_P, m_0}$ . The value of  $d_{P, k}$  at a query point  $x$  is given by

$$d_{P, k}^2(x) = \frac{1}{k} \sum_{p \in \text{NN}_P^k(x)} \|x - p\|^2,$$

where  $\text{NN}_P^k(x) \subseteq P$  denotes the  $k$  nearest neighbors in  $P$  to the point  $x \in \mathbb{R}^d$ . (Note that while the  $k$ -th nearest neighbor itself might be ambiguous, on the boundary of an order- $k$  Voronoi cell, the distance to the  $k$ -th nearest neighbor is always well defined, and so is  $d_{P, k}$ .)

The most important property of the distance function  $d_{\mu, m_0}$  is its stability, for a fixed  $m_0$ , under perturbations of the underlying measure  $\mu$ . This property provides a bridge between the underlying (continuous)  $\mu$  and the discrete measure  $\mathbf{1}_P$ . According to [6, Theorem 3.5], for any two probability measures  $\mu$  and  $\nu$  on  $\mathbb{R}^d$ ,

$$(2) \quad \|d_{\mu, m_0} - d_{\nu, m_0}\|_{\infty} \leq m_0^{-1/2} W_2(\mu, \nu),$$

where  $W_2(\mu, \nu)$  denotes the Wasserstein distance between the two measures. The bound in this inequality depends on the choice of  $m_0$ , which acts as a smoothing parameter.

### 3. WITNESSED $k$ -DISTANCE

In this section, we describe a simple scheme for approximating the distance to a uniform measure, together with a general error bound. The main contribution of our work, presented in Section 4, is the analysis of the quality of approximation given by this scheme when the input points come from a measure concentrated on a lower-dimensional subset of the Euclidean space.

**3.1.  $k$ -Distance as a Power Distance.** Given a set of points  $U = \{u_1, \dots, u_n\}$  in  $\mathbb{R}^d$  with weights  $(w_u)$ , we call *power distance* to  $U$  the function  $\text{pow}_U$  obtained as the lower envelope of all the functions  $d_u : x \mapsto \|u - x\|^2 - w_u$ , where  $u$  ranges over  $U$ . By Proposition 3.1 in [6], we can express the square of any distance to a measure as a power distance with non-positive weights. The following proposition recalls this property in the case of  $k$ -distance; it's equivalent to the well-known fact that the order- $k$  Voronoi diagrams can be written as the power diagrams for a certain set of points and weights [3].

**Proposition 1.** *For any  $P \subseteq \mathbb{R}^d$ , denote by  $\text{Bary}^k(P)$  the set of barycenters of any subset of  $k$  points in  $P$ . Then*

$$(3) \quad d_{P, k}^2 = \min \left\{ \|x - \bar{p}\|^2 - w_{\bar{p}}; \bar{p} \in \text{Bary}^k(P) \right\},$$

where the weight of a barycenter  $\bar{p} = \frac{1}{k} \sum_i p_i$  is given by  $w_{\bar{p}} := -\frac{1}{k} \sum_i \|\bar{p} - p_i\|^2$ .

*Proof.* Given  $k$  distinct data points  $p_1, \dots, p_k$  in  $P$ , denote their isobarycenter by  $\bar{p}$  and consider the function  $d_{\bar{p}}^2(x) := \frac{1}{k} \sum_{1 \leq i \leq k} \|x - p_i\|^2$ . An easy computation shows that

$$d_{\bar{p}}^2(x) = \frac{1}{k} \sum_{1 \leq i \leq k} \|x - p_i\|^2 = \|x - \bar{p}\|^2 - w_{\bar{p}},$$

where the weight  $w_{\bar{p}} = -\frac{1}{k} \sum_{1 \leq i \leq k} \|\bar{p} - p_i\|^2$ . The proposition follows since  $d_{\bar{p},k}^2$  is the minimum of all the functions  $d_{\bar{p}}$ .  $\square$

In other words, the square of the  $k$ -distance function to  $P$  coincides exactly with the power distance to the set of barycenters  $\text{Bary}^k(P)$  with the weights defined above. From this expression, it follows that the sublevel sets of the  $k$ -distance  $d_{P,k}$  are finite unions of balls,

$$d_{P,k}^{-1}([0, \rho]) = \bigcup_{\{p_i\} \in \text{NN}_P^k(\mathbb{R}^d)} \text{B}(\bar{p}, (\rho^2 + w_{\bar{p}})^{1/2}).$$

Therefore, ignoring the complexity issues, it is possible to compute the homotopy type of this sublevel set by considering the weighted alpha-shape of  $\text{Bary}^k(P)$  [15], which is a subcomplex of the regular triangulation of the set of weighted barycenters.

From the proof of Proposition 1, we also see that the only barycenters that actually play a role in Equation (3) are the barycenters of  $k$  points of  $P$  whose order- $k$  Voronoi cell is not empty. However, the dependence on the number of non-empty order- $k$  Voronoi cells makes computation intractable even for moderately sized point clouds in the Euclidean space [9]. One way to avoid this difficulty is to replace the  $k$ -distance to  $P$  by an approximate  $k$ -distance, defined as in Equation (3), but where the minimum is taken over a smaller set of barycenters. Then, the question is: Given a point set  $P$ , can we replace the set of barycenters  $\text{Bary}_P^k$  in the definition of  $k$ -distance by a small subset  $B$  while controlling the approximation error  $\|\text{pow}_B^{1/2} - d_{P,k}\|_\infty$ ?

Replacing the  $k$ -distance with another power distance is especially attractive since many geometric and topological inference methods relying on distance functions continue to hold when one of the functions is replaced by a good approximation *in the class of power distances*. When this is the case, and some sampling conditions are met, it is possible, for instance, to recover the homotopy type of the underlying compact set (see the Reconstruction Theorem in [6].)

**3.2. Approximating by witnessed  $k$ -distance.** We consider a subset of the barycenters suggested by the input data. The answer to our question is affirmative when the input point cloud  $P$  satisfies the hypotheses **(H1)**-**(H3)**.

**DEFINITION 2.** For every point  $p$  in  $P$ , the barycenter of  $p$  and its  $(k-1)$  nearest neighbors in  $P$  is called a *witnessed  $k$ -barycenter*. Let  $\text{Bary}_w^k(P)$  be the set of all such barycenters. We get one witnessed barycenter for every point of the sampled point set, and define the *witnessed  $k$ -distance*,

$$d_{P,k}^w = \min\{\|x - \bar{p}\|^2 - w_{\bar{p}}; \bar{p} \in \text{Bary}_w^k(P)\}.$$

Computing the set of all witnessed barycenters of a point set  $P$  requires only finding the  $k-1$  nearest neighbors of every point in  $P$ . This search problem has

a long history in computational geometry [2, 8, 18], and now has several practical implementations.

**General error bound.** Because the distance functions we consider are defined by minima, and  $\text{Bary}_w^k(P)$  is a subset of  $\text{Bary}^k(P)$ , the witnessed  $k$ -distance is always greater than the exact  $k$ -distance. In the lemma below, we give a general multiplicative upper bound. This lemma does not assume any special property for the input point set  $P$ . However, even such a coarse bound can be used to estimate Betti numbers of sublevel sets of  $d_{P,k}$ , using arguments similar to those in [7].

**Lemma 1** (General Bound<sup>1</sup>). *For any finite point set  $P \subseteq \mathbb{R}^d$  and  $0 < k < |P|$ ,*

$$0 \leq d_{P,k}^w - d_{P,k} \leq 2d_P \leq 2d_{P,k}.$$

*Proof.* Let  $y$  be a point in  $\mathbb{R}^d$ , and  $\bar{p}$  the barycenter associated to an order- $k$  Voronoi cell containing  $y$ , i.e.  $\bar{p}$  is such that  $d_{P,k}(y) = d_{\bar{p}}(y) := (\|x - \bar{p}\|^2 - w_{\bar{p}})^{1/2}$ . Let us find a witnessed barycenter  $\bar{q}$  close to  $\bar{p}$ . By definition,  $\bar{p}$  is the barycenter of the  $k$  nearest neighbors  $p_1, \dots, p_k$  of  $y$  in  $P$ . Let  $x := p_1$  be the nearest neighbor of  $y$ , and  $\bar{q}$  the barycenter witnessed by the point  $x$ . Then,  $d_P(y) = \|x - y\| \leq d_{P,k}(y)$ , and

$$\begin{aligned} d_{P,k}^w(y) &\leq d_{\bar{q}}(y) \leq d_{\bar{q}}(x) + \|x - y\| \leq d_{\bar{p}}(x) + \|x - y\| \\ &\leq d_{\bar{p}}(y) + 2\|x - y\| = d_{P,k}(y) + 2d_P(y). \end{aligned}$$

Since  $y$  was chosen arbitrarily, the claim follows.  $\square$

#### 4. APPROXIMATION QUALITY

Let us recall briefly our hypothesis **(H1)**-**(H3)**. There is an ideal, well-conditioned measure  $\mu$  on  $\mathbb{R}^d$  supported on an unknown compact set  $K$ . We also have a noisy version of  $\mu$ , that is another measure  $\nu$  with  $W_2(\mu, \nu) \leq \sigma$ , and we suppose that our data set  $P$  consists of  $N$  points independently sampled from  $\nu$ . In this section we give conditions under which the witnessed  $k$ -distance to  $P$  provides a good approximation of the distance to the underlying set  $K$ .

**4.1. Dimension of a measure.** First, we make precise the main assumption **(H1)** on the underlying measure  $\mu$ , which we use to bound the approximation error made when replacing the exact by the witnessed  $k$ -distance. We require  $\mu$  to be low dimensional in the following sense.

**DEFINITION 3.** A measure  $\mu$  on  $\mathbb{R}^d$  is said to have *dimension at most  $\ell$* , which we denote by  $\dim \mu \leq \ell$ , if there is a positive constant  $\alpha_\mu$  such that the amount of mass contained in the ball  $B(p, r)$  is at least  $\alpha_\mu r^\ell$ , for every point  $p$  in the support of  $\mu$  and every radius  $r$  smaller than the diameter of this support.

The important assumption here is that the lower bound  $\mu(B(p, r)) \geq \alpha r^\ell$  should be true for some positive constant  $\alpha$  and for  $r$  smaller than a given constant  $R$ . The choice of  $R = \text{diam}(\text{spt}(\mu))$  provides a normalization of the constant  $\alpha_\mu$  and slightly simplifies the statements of the results.

Let  $M$  be an  $\ell$ -dimensional compact submanifold of  $\mathbb{R}^d$ , and  $f : M \rightarrow \mathbb{R}$  a positive weight function on  $M$  with values bounded away from zero and infinity. Then, the dimension of the volume measure on  $M$  weighted by the function  $f$  is at most  $\ell$ .

<sup>1</sup>The authors thank Daniel Chen for strengthening the earlier version of this bound.

A quantitative statement can be obtained using the Bishop-Günther comparison theorem; the bound depends on the maximum absolute sectional curvature of the manifold  $M$ , as shown in Proposition 4.9 in [6]. Note that the positive lower bound on the density is really necessary. For instance, the dimension of the standard Gaussian distribution  $\mathcal{N}(0, 1)$  on the real line is not bounded by 1, nor by any positive constant, because the density of this distribution decreases to zero faster than any function  $r \mapsto 1/r^\ell$  as one moves away from the origin.

It is easy to see that if  $m$  measures  $\mu_1, \dots, \mu_m$  have dimension at most  $\ell$ , then so does their sum. Consequently, if  $(M_j)$  is a finite family of compact submanifolds of  $\mathbb{R}^d$  with dimensions  $(d_j)$ , and  $\mu_j$  is the volume measure on  $M_j$  weighted by a function bounded away from zero and infinity, the dimension of the measure  $\mu = \sum_{j=1}^m \mu_j$  is at most  $\max_j d_j$ .

**4.2. Bounds.** In the remainder of this section, we bound the error between the witnessed  $k$ -distance  $d_{P,k}^w$  and the (ordinary) distance  $d_K$  to the compact set  $K$ . We start from a proposition from [6] that bounds the error between the exact  $k$ -distance  $d_{P,k}$  and  $d_K$ :

**Theorem 1** (Exact Bound). *Let  $\mu$  denote a probability measure with dimension at most  $\ell$ , and supported on a set  $K$ . Consider the uniform measure  $\mathbf{1}_P$  on a point cloud  $P$ , and set  $m_0 = k/|P|$ . Then*

$$\|d_{P,k} - d_K\|_\infty \leq m_0^{-1/2} W_2(\mu, \mathbf{1}_P) + \alpha_\mu^{-1/\ell} m_0^{1/\ell}.$$

*Proof.* Recall that  $d_{P,k} = d_{\mathbf{1}_P, m_0}$ . Using the triangle inequality and Equation (2), one has

$$\begin{aligned} \|d_{\mathbf{1}_P, m_0} - d_K\|_\infty &\leq \|d_{\mu, m_0} - d_{\mathbf{1}_P, m_0}\|_\infty + \|d_{\mu, m_0} - d_K\|_\infty \\ &\leq m_0^{-1/2} W_2(\mu, \mathbf{1}_P) + \|d_{\mu, m_0} - d_K\|_\infty \end{aligned}$$

Then, from Lemma 4.7 in [6],  $\|d_{\mu, m_0} - d_K\|_\infty \leq \alpha_\mu^{-1/\ell} m_0^{1/\ell}$ , and the claim follows.  $\square$

To make this bound concrete, let us construct a simple example where the term corresponding to the Wasserstein noise and the term corresponding to the smoothing have the same order of magnitude.

**EXAMPLE.** Consider the restriction  $\mu$  of the Lebesgue measure to the  $\ell$ -dimensional unit ball  $K := B(0, 1)$ , rescaled to be a probability measure by a factor  $1/\text{vol}^\ell B(0, 1)$ . For a given mass parameter  $m_0$ , consider the second measure  $\nu$  obtained by moving every bit of mass of  $\mu$  in the  $\ell$ -ball  $B(0, m_0^{1/\ell})$  to the closest point in the  $(\ell-1)$ -sphere  $S(0, m_0^{1/\ell})$ . By construction,

$$\begin{aligned}
W_2(\mu, \nu)^2 &= \int_{B(0, m_0^{1/\ell})} (m_0^{1/\ell} - \|x\|)^2 d\mu(x) \\
&= \frac{\text{vol}^\ell B(0, 1)}{\text{vol}^{\ell-1} S(0, 1)} \int_0^{m_0^{1/\ell}} r^{\ell-1} (m_0^{1/\ell} - r)^2 dr \\
&= \ell \int_0^{m_0^{1/\ell}} (r^{\ell+1} - 2r^\ell m_0^{1/\ell} + r^{\ell-1} m_0^{2/\ell}) dr \\
&= \frac{2m_0^{1+2/\ell}}{(\ell+1)(\ell+2)}
\end{aligned}$$

The distance  $d_{\nu, m_0}(0)$  of the origin to  $\nu$  is easy to compute: the radius of the smallest ball centered at the origin with a mass  $m_0$  of  $\nu$  is exactly  $m_0^{1/\ell}$ . Hence,

$$\begin{aligned}
\|d_K - d_{\nu, m_0}\|_\infty &\geq |d_K(0) - d_{\nu, m_0}(0)| \\
&= m_0^{1/\ell} = C_\ell m_0^{-1/2} W_2(\mu, \nu).
\end{aligned}$$

In other words, the two terms in the bound in Theorem 1 differ by a constant factor.  $\square$

In the main theorem of this section, the exact  $k$ -distance in Theorem 1 is replaced by the witnessed  $k$ -distance. Observe that the new error term is only a constant factor off from the old one.

**Theorem 2** (Witnessed Bound). *Let  $\mu$  be a probability measure satisfying the dimension assumption and let  $K$  be its support. Consider the uniform measure  $\mathbf{1}_P$  on a point cloud  $P$ , and set  $m_0 = k/|P|$ . Then,*

$$\|d_{P, k}^w - d_K\|_\infty \leq 3m_0^{-1/2} W_2(\mu, \mathbf{1}_P) + 12m_0^{1/\ell} \alpha_\mu^{-1/\ell}.$$

Before proving the theorem, we start with an auxiliary lemma showing that a measure  $\nu$  close to a measure  $\mu$  satisfying an upper dimension bound (as in Definition 3) remains concentrated around the support of  $\mu$ .

**Lemma 2** (Concentration). *Let  $\mu$  be a probability measure satisfying the dimension assumption, and  $\nu$  be another probability measure. Let  $m_0$  be a mass parameter. Then, for every point  $p$  in the support of  $\mu$ ,  $\nu(B(p, \eta)) \geq m_0$ , where  $\eta = m_0^{-1/2} W_2(\mu, \nu) + 4m_0^{1/2+1/\ell} \alpha_\mu^{-1/\ell}$ .*

*Proof.* Let  $\pi$  be an optimal transport plan between  $\nu$  and  $\mu$ . For a fixed point  $p$  in the support  $K$  of  $\mu$ , let  $r$  be the smallest radius such that  $B(p, r)$  contains at least  $2m_0$  of mass  $\mu$ . Consider now a submeasure  $\mu'$  of  $\mu$  of mass exactly  $2m_0$  and whose support is contained in the ball  $B(p, r)$ . This measure is obtained by transporting a submeasure  $\nu'$  of  $\nu$  by the optimal transport plan  $\pi$ . Our goal is to determine for what choice of  $\eta$  the ball  $B(p, \eta)$  contains a  $\nu'$ -mass (and, therefore, a  $\nu$ -mass) of at least  $m_0$ . We make use of the Chebyshev's inequality for  $\nu'$  to bound the mass of  $\nu'$  outside of the ball  $B(p, \eta)$ :

$$\begin{aligned}
\nu'(\mathbb{R}^d \setminus B(p, \eta)) &= \nu'(\{x \in \mathbb{R}^d; \|x - p\| \geq \eta\}) \\
(4) \quad &\leq \frac{1}{\eta^2} \int \|x - p\|^2 d\nu'
\end{aligned}$$

Observe that the right hand term of this inequality is exactly the Wasserstein distance between  $\nu'$  and the Dirac mass  $2m_0\delta_p$  divided by  $\eta^2$ . We bound this Wasserstein distance using the triangle inequality:

$$(5) \quad \begin{aligned} \int \|x - p\|^2 d\nu' &= W_2^2(\nu', 2m_0\delta_p) \\ &\leq (W_2(\mu', \nu') + W_2(\mu', 2m_0\delta_p))^2 \\ &\leq (W_2(\mu, \nu) + 2m_0r)^2 \end{aligned}$$

Combining Equation (4) and Equation (5), we get:

$$\begin{aligned} \nu(\bar{B}(p, \eta)) &\geq \nu'(\bar{B}(p, \eta)) \geq \nu'(\mathbb{R}^d) - \nu'(\mathbb{R}^d \setminus B(p, \eta)) \\ &\geq 2m_0 - \frac{(W_2(\mu, \nu) + 2m_0r)^2}{\eta^2}. \end{aligned}$$

By the lower bound on the dimension of  $\mu$ , and the definition of the radius  $r$ , one has  $r \leq (2m_0/\alpha_\mu)^{1/\ell}$ . Hence, the ball  $\bar{B}(p, \eta)$  contains a mass of at least  $m_0$  as soon as

$$\frac{(W_2(\mu, \nu) + \alpha_\mu^{-1/\ell} 2^{1+1/\ell} m_0^{1+1/\ell})^2}{\eta^2} \leq m_0.$$

This will be true, in particular, if  $\eta$  is larger than

$$W_2(\mu, \nu)m_0^{-1/2} + 4\alpha_\mu^{-1/\ell} m_0^{1/2+1/\ell}. \quad \square$$

*Proof of the Witnessed Bound Theorem.* Since the witnessed  $k$ -distance is a minimum over fewer barycenters, it is larger than the real  $k$ -distance. Using this fact and the Exact Bound Theorem one gets the lower bound:

$$d_{P,k}^w \geq d_{P,k} \geq d_K - \left( m_0^{-1/2} W_2(\mu, \mathbf{1}_P) + \alpha_\mu^{-1/\ell} m_0^{1/\ell} \right)$$

For the upper bound, choose  $\eta$  as given by the previous Lemma 2. Then, for every point  $y$  in  $K$ , the ball  $B(p, \eta)$  contains at least  $k$  points in the point cloud  $P$ . Let  $p_1$  be one of these points, and  $p_2, \dots, p_k$  be the  $(k-1)$  nearest neighbors of  $p_1$  in  $P$ . The points  $p_2, \dots, p_k$  cannot be at a distance greater than  $2\eta$  from  $p_1$ , and, consequently, cannot be at a distance greater than  $3\eta$  from  $y$ . By definition, the barycenter  $\bar{p}$  of the points  $\{p_i\}$  is witnessed by  $p_1$ . Hence,

$$d_{P,k}^w(y) \leq d_{\bar{p}}(y) := \left( \frac{1}{k} \sum_{i=1}^k \|y - p_i\|^2 \right)^{1/2} \leq 3\eta$$

Since  $d_{P,k}^w$  is 1-Lipschitz, we get  $d_{P,k}^w(x) \leq 3\eta + \|y - x\|$ . This inequality is true for every point  $y$  in  $K$ ; minimizing over all such  $y$ , we obtain  $d_{P,k}^w(x) \leq 3\eta + d_K(x)$ . Recall that  $m_0 \leq 1$ , as is  $m_0^{1/2}$ . To match the form of the bound in Theorem 1, we drop  $m_0^{1/2}$  from the second term of  $\eta$  in the Concentration Lemma. Substituting the result into the last inequality, we complete the proof.  $\square$

## 5. CONVERGENCE UNDER EMPIRICAL SAMPLING

One term remains moot in the bounds in Theorems 1 and 2, namely the Wasserstein distance  $W_2(\mu, \mathbf{1}_P)$ . In this section, we analyze its convergence. The rate depends on the complexity of the measure  $\mu$ , defined below. The moral of this

section is that if a measure can be well approximated with few points, then it is also well approximated by random sampling.

**DEFINITION 4.** The *complexity* of a probability measure  $\mu$  at a scale  $\varepsilon > 0$  is the minimum cardinality of a finitely supported probability measure  $\nu$  which  $\varepsilon$ -approximates  $\mu$  in the Wasserstein sense, i.e. such that  $W_2(\mu, \nu) \leq \varepsilon$ . We denote this number by  $\mathcal{N}_\mu(\varepsilon)$ .

Observe that this notion is very close to the  $\varepsilon$ -covering number of a compact set  $K$ , denoted by  $\mathcal{N}_K(\varepsilon)$ , which counts the minimum number of balls of radius  $\varepsilon$  needed to cover  $K$ . It's worth noting that if measures  $\mu$  and  $\nu$  are close — as are the measure  $\mu$  and its noisy approximation  $\nu$  in the previous section — and  $\mu$  has low complexity, then so does the measure  $\nu$ . The following lemma shows that measures satisfying the dimension assumption have low complexity. Its proof follows from a classical covering argument that appears, for example, in Proposition 4.1 of [19].

**Lemma 3** (Dimension–Complexity). *Let  $K$  be the support of a measure  $\mu$  with  $\dim \mu \leq \ell$ . Then,*

- (i) *for every positive  $\varepsilon$ ,  $\mathcal{N}_K(\varepsilon) \leq \alpha_\mu/\varepsilon^\ell$ . Said otherwise, the upper box-counting dimension of  $K$  is bounded:*

$$\dim(K) := \limsup_{\varepsilon \rightarrow 0} \log(\mathcal{N}_K(\varepsilon)) / \log(1/\varepsilon) \leq \ell.$$

- (ii) *for every positive  $\varepsilon$ ,  $\mathcal{N}_\mu(\varepsilon) \leq \alpha_\mu 5^\ell / \varepsilon^\ell$ .*

**Theorem 3** (Convergence). *Let  $\mu$  be a probability measure on  $\mathbb{R}^d$  whose support has diameter at most  $D$ , and let  $P$  be a set of  $N$  points independently drawn from the measure  $\mu$ . Then, for  $\varepsilon > 0$ ,*

$$\mathbb{P}(W_2(\mathbf{1}_P, \mu) \leq 4\varepsilon) \geq 1 - \mathcal{N}_\mu(\varepsilon) \exp(-2N\varepsilon^2/(D\mathcal{N}_\mu(\varepsilon))^2) - \exp(-2N\varepsilon^4/D^2).$$

*Proof.* Let  $n$  be a fixed integer, and let  $\varepsilon$  be the minimum Wasserstein distance between  $\mu$  and a measure  $\bar{\mu}$  supported on (at most)  $n$  points. Let  $S$  be the support of the optimal measure  $\bar{\mu}$ , so that  $\bar{\mu}$  can be decomposed as  $\sum_{s \in S} \alpha_s \delta_s$  ( $\alpha_s \geq 0$ ). Let  $\pi$  be an optimal transport plan between  $\mu$  and  $\bar{\mu}$ ; this is equivalent to finding a decomposition of  $\mu$  as a sum of  $n$  non-negative measures  $(\pi_s)_{s \in S}$  such that  $\text{mass}(\pi_s) = \alpha_s$ , and

$$\sum_{s \in S} \int \|x - s\|^2 d\pi_s(x) = \varepsilon^2 = W_2(\mu, \bar{\mu})^2.$$

Drawing a random point  $X$  from the measure  $\mu$  amounts to (i) choosing a random point  $s$  in the set  $S$  (with probability  $\alpha_s$ ) and (ii) drawing a random point  $X$  following the distribution  $\pi_s$ . Given  $N$  independent points  $X_1, \dots, X_N$  drawn from the measure  $\mu$ , denote by  $I_{s,N}$  the proportion of the  $(X_i)$  for which the point  $s$  was selected in step (i). Hoeffding's inequality allows to easily quantify how far the proportion  $I_{s,N}$  deviates from  $\alpha_s$ :  $\mathbb{P}(|I_{s,N} - \alpha_s| \geq \delta) \leq \exp(-2N\delta^2)$ . Combining these inequalities for every point  $s$  and using the union bound yields

$$\mathbb{P}\left(\sum_{s \in S} |I_{s,N} - \alpha_s| \leq \delta\right) \geq 1 - n \exp(-2N\delta^2/n^2).$$

For every point  $s$ , denote by  $\tilde{\pi}_s$  the distribution of the distances to  $s$  in the submeasure  $\pi_s$ , i.e. the measure on the real line defined by  $\tilde{\pi}_s(I) := \pi_s(\{x \in$

$\mathbb{R}^d; \|x - s\| \in I\}$ ) for every interval  $I$ . Define  $\tilde{\mu}$  as the sum of the  $\tilde{\pi}_s$ ; by the change of variable formula one has

$$\int_{\mathbb{R}^d} t^2 d\tilde{\mu}(t) = \sum_s \int_{\mathbb{R}^d} t^2 d\tilde{\pi}_s = \sum_s \int_{\mathbb{R}^d} \|x - s\|^2 d\pi_s = \varepsilon^2.$$

Given a random point  $X_i$  sampled from  $\mu$ , denote by  $Y_i$  Euclidean distance between the point  $X_i$  and the point  $s$  chosen in step (i). By construction, the distribution of  $Y_i$  is given by the measure  $\tilde{\mu}$ ; using the Hoeffding inequality again one gets

$$\mathbb{P}\left(\frac{1}{N} \sum_{i=1}^N Y_i^2 \geq (\varepsilon + \eta)^2\right) \leq 1 - \exp(-2N\eta^2\varepsilon^2/D^2).$$

In order to conclude, we need to define a transport plan from the empirical measure  $\mathbf{1}_P = \frac{1}{N} \sum_{i=1}^N \delta_{X_i}$  to the finite measure  $\bar{\mu}$ . To achieve this, we order the points  $(X_i)$  by increasing distance  $Y_i$ ; then transport every Dirac mass  $\frac{1}{N} \delta_{X_i}$  to the corresponding point  $s$  in  $S$  until  $s$  is “full”, i.e. the mass  $\alpha_s$  is reached. The squared cost of this transport operation is at most  $\frac{1}{N} \sum_{i=1}^N Y_i^2$ . Then distribute the remaining mass among the  $s$  points in any way; the cost of this step is at most  $D$  times  $\sum_{s \in S} |I_{s,N} - \alpha_s|$ . The total cost of this transport plan is the sum of these two costs. From what we have shown above, setting  $\eta = \varepsilon$  and  $\delta = \varepsilon/D$ , one gets

$$\mathbb{P}(W_2(\mathbf{1}_P, \mu) \leq 4\varepsilon) \geq 1 - n \exp(-2N\varepsilon^2/(Dn)^2) - \exp(-2N\varepsilon^4/D^2). \quad \square$$

As a consequence of the Dimension–Complexity Lemma 3 and of the Convergence Theorem 3, any measure  $\mu$  satisfying an upper bound on its dimension is well approximated by empirical sampling. A result similar to the Convergence Theorem follows when the samples are drawn not from the original measure  $\mu$ , but from a “noisy” approximation  $\nu$  which need not be compactly supported.

**DEFINITION 5.** A probability measure  $\nu$  on  $\mathbb{R}^d$  is a  $\sigma$ -*perturbation* of another probability measure  $\mu$  if there exist two random vectors  $X$  and  $Y$  whose distributions are respectively  $\mu$  and  $\nu$  such that:

- (i) The expectation  $\mathbb{E}(\|X - Y\|)^2$  is at most  $\sigma^2$ . This implies, in particular, that the Wasserstein distance between  $\mu$  and  $\nu$  is at most  $\sigma$ .
- (ii) The random variable  $Z := \|X - Y\|^2$  is sub-exponential,  $\mathbb{P}(Z \geq x) \leq \exp(-x/\sigma)$ .

**EXAMPLE.** Consider the measure  $\nu$  obtained by convolving  $\mu$  with an isotropic Gaussian distribution of variance  $\sigma^2$ . Consider the random variables  $X$  and  $Y$  associated to the transport plan given by the convolution. Then,  $\mathbb{E}(\|X - Y\|^2) \leq \sigma^2$ . Moreover, we have  $\mathbb{P}(\|X - Y\| \geq x) = \text{erfc}(x/\sigma) \leq \exp(-x^2/\sigma^2)$ , where  $\text{erfc}$  is the complementary error function. This yields  $\mathbb{P}(Z \geq x) \leq \exp(-x/\sigma)$ , i.e. the probability measure  $\nu$  is a  $\sigma$ -perturbation of  $\mu$ .

**Corollary 1** (Noisy Convergence). *Let measure  $\nu$  be a  $\sigma$ -perturbation of a measure  $\mu$  satisfying the conditions of Theorem 3. Let  $Q$  be a set of  $N$  points drawn independently from the measure  $\nu$ . Then,*

$$\mathbb{P}(W_2(\mathbf{1}_Q, \mu) \leq 6\sigma) \geq 1 - \exp(-N\sigma) - \mathcal{N}_\mu(\sigma) \exp(-2N\sigma^2/(DN_\mu(\sigma))^2) - \exp(-2N\sigma^4/D^2).$$

*Proof.* Consider the two sequences of iid. random variables  $(X_i, Y_i)$  obtained by independent (in  $i$ ) copies of the two random variables  $(X, Y)$ , with distribution  $\mu$  and  $\nu$  introduced in Definition 5. The main difficulty is bounding the probability of the Wasserstein distance between the uniform probability measures on  $P = \{X_1, \dots, X_n\}$  and  $Q = \{Y_1, \dots, Y_n\}$  exceeding  $(1 + \varepsilon)\sigma$ . We apply the standard Chernoff technique, based on moment generating functions (for example, see Problem 1.8 in [14]). Using the sub-exponential property, we can bound the moment generating function of the random variable  $Z_i = \|X_i - Y_i\|^2$  for  $t < 1/\sigma$ :

$$\begin{aligned} \mathbb{E}(e^{tZ_i}) &= \int_{\mathbb{R}^+} \left( \frac{d}{dx} \exp(tx) \right) \mathbb{P}(Z_i \geq x) dx \\ &\leq \int_{\mathbb{R}^+} t \exp(tx) \exp(-x/\sigma) dx = \left[ t \frac{\exp((t - 1/\sigma)x)}{t - 1/\sigma} \right]_0^\infty = \frac{t}{1/\sigma - t}. \end{aligned}$$

This yields the following tail inequality for  $\bar{Z} = Z_1 + \dots + Z_N$ :

$$\mathbb{P}(\bar{Z} \geq x) \leq \frac{\mathbb{E}(e^{t\bar{Z}})}{e^{tx}} = \exp \left[ N \ln \left( \frac{\sigma t}{1 - \sigma t} \right) - tx \right].$$

With  $x = N(1 + \varepsilon)^2 \sigma^2$  and  $t = 1/(2\sigma)$  we obtain

$$\begin{aligned} \mathbb{P} \left( \frac{1}{N} \bar{Z} \geq (1 + \varepsilon)^2 \sigma^2 \right) &\leq \exp \left[ N \ln \left( \frac{\sigma t}{1 - \sigma t} \right) - tN(1 + \varepsilon)^2 \sigma^2 \right] \\ &= \exp \left( -\frac{1}{2} N(1 + \varepsilon)^2 \sigma \right). \end{aligned}$$

As a consequence, setting  $\varepsilon = 1$ , we obtain  $\mathbb{P}(W_2(\mathbf{1}_P, \mathbf{1}_Q) \geq 2\sigma) \leq \exp(-N\sigma)$ . Applying Theorem 3, we get:

$$\begin{aligned} \mathbb{P}(W_2(\mathbf{1}_Q, \mu) \leq 6\sigma) &\geq \mathbb{P}(W_2(\mathbf{1}_P, \mathbf{1}_Q) + W_2(\mathbf{1}_P, \mu) \leq 6\sigma) \\ &\geq \mathbb{P}(W_2(\mathbf{1}_P, \mathbf{1}_Q) < 2\sigma) \cdot \mathbb{P}(W_2(\mathbf{1}_P, \mu) \leq 4\sigma) \\ &\geq (1 - \exp(-N\sigma)) \cdot (1 - \mathcal{N}_\mu(\sigma) \exp(-2N\sigma^2/(D\mathcal{N}_\mu(\sigma))^2) - \exp(-2N\sigma^4/D^2)) \\ &\geq 1 - \exp(-N\sigma) - \mathcal{N}_\mu(\sigma) \exp(-2N\sigma^2/(D\mathcal{N}_\mu(\sigma))^2) - \exp(-2N\sigma^4/D^2). \end{aligned}$$

□

We combine Theorem 2 (Witnessed Bound), Corollary 1 (Noisy Convergence) and Lemma 3 (Dimension-Complexity) to get the following probabilistic statement.

**Theorem 4** (Approximation). *Suppose that  $\mu$  is a measure satisfying the dimension assumption, supported on a set  $K$  of diameter  $D$ , and  $\nu$  is a  $\sigma$ -perturbation of  $\mu$ . Let  $P$  be a set of  $N$  points independently sampled from  $\nu$ . Then, the inequality*

$$\|d_{P,k}^w - d_K\|_\infty \leq 18m_0^{-1/2} \sigma + 12m_0^{1/\ell} \alpha_\mu^{-1/\ell}$$

*holds with probability at least*

$$1 - (\alpha_\mu 5^\ell / \sigma^\ell) \exp(-2N\sigma^{2+2\ell}/(D\alpha_\mu 5^\ell)^2) - \exp(-2N\sigma^4/D^2) - \exp(-N\sigma).$$

*Proof.* Lemma 3 upper-bounds the covering number  $\mathcal{N}_\mu(\sigma)$  by  $\alpha_\mu 5^\ell / \sigma^\ell$ . Substituting it into the bound in Corollary 1 gives the probabilistic bound on the Wasserstein distance of  $6\sigma$  between the measures  $\mathbf{1}_P$  and  $\mu$ . Plugging it into Theorem 2 yields the claimed result. □

## 6. DISCUSSION

We illustrate the utility of the Witnessed Bound Theorem by example and an inference statement. Figure 1(a) shows 6000 points drawn from the uniform distribution on a sideways figure-8 (in red), convolved with a Gaussian distribution. The ordinary distance function has no hope of recovering geometric information out of these points since both loops of the figure-8 are filled in. In Figure 1(b), we show the sublevel sets of the distance to the uniform measure on the point set, both the witnessed  $k$ -distance and the exact  $k$ -distance. Both functions recover the topology of figure-8; the bits missing from the witnessed  $k$ -distance smooth out the boundary of the sublevel set, but do not affect the image at large.

**Complexes.** Since we are working with the power distance to a weighted point set (the witnessed barycenters  $U$ ), we can employ different simplicial complexes commonly used in the computational geometry literature [17]. Recall that an (abstract) simplex is a subset of some universal set, in our case the witnessed barycenters  $U$ ; a simplicial complex is a collection of simplices, closed under the face relation.

The simplest construction is the Čech complex. It contains a simplex if the balls defined by the points at the power distance  $r$  from the witnessed barycenters intersect:

$$\check{C}_r(U) = \{\sigma \subseteq U \mid \bigcap_{u \in \sigma} B(u, (r^2 + w_u)^{1/2})\}.$$

A closely related geometric construction, the weighted alpha complex, is defined by clipping these balls using the power diagram of the witnessed barycenters, see [15]. By the Nerve Theorem [17], both the Čech complex  $\check{C}_r(U)$  and the alpha complex are homotopy equivalent to the sublevel sets of the power distance to  $U$ ,  $\text{pow}_U^{-1}(-\infty, r]$ .

In many applications, points are only given through their pairwise distances, rather than explicit coordinates. For this reason and because of its computational simplicity, the Vietoris-Rips complex is a popular choice. This complex is defined as the flag (or clique) complex of the 1-skeleton of the Čech complex. Said otherwise, a simplex  $\sigma$  belongs to the Vietoris-Rips complex iff all its edges belong to the Čech complex, i.e.

$$\text{VR}_r(U) = \{\sigma \subseteq U \mid \{u, v\} \in \check{C}_r(U) \text{ for all } u, v \in \sigma\}.$$

In the case of the witnessed  $k$ -distance, the pairwise distances between the input points suffice for the construction of the Vietoris-Rips complex on the witnessed barycenters; we give the details in the Appendix A.

Note that the Vietoris-Rips complex  $\text{VR}_r(U)$  does not, in general, have the homotopy type of  $\text{pow}_U^{-1}(-\infty, r]$ . It is, however, possible to prove inference results for homology given an *interleaving property*, i.e. there exists a constant  $\alpha \geq 1$  such that  $\check{C}_r(U) \subseteq \text{VR}_r(U) \subseteq \check{C}_{\alpha r}(U)$ . The inclusion  $\check{C}_r(U) \subseteq \text{VR}_r(U)$  always holds, simply by definition. However, the second inclusion does not necessary hold if the weights are positive, as the following example demonstrates.

**EXAMPLE.** Consider the weighted point set  $U$  made of the three vertices  $(u, v, w)$  of an equilateral triangle with unit sidelength and weights  $w_u = w_v = w_w = 1/4$ . Then, for any non-negative  $r$ , the Vietoris-Rips complex  $\text{VR}_r(U)$  contains the triangle, while the Čech complex  $\check{C}_r(U)$  contains this triangle only as soon as

$(r^2 + 1/4)^{1/2} \geq 1/\sqrt{3}$ , i.e.  $r \geq 1/\sqrt{12}$ . In this case, there is no  $\alpha$  such that the inclusion  $\text{VR}_r(U) \subseteq \check{\text{Cech}}_{\alpha r}(U)$  holds for every positive  $r$ .

On the other hand, the following lemma shows that when the weights  $(w_u)_{u \in U}$  are non-positive, the inclusion  $\text{VR}_r(U) \subseteq \check{\text{C}}_{2r}(U)$  always holds. This property lets us extend the usual homology inference results from Vietoris-Rips complexes to the (weighted) Vietoris-Rips complexes associated to the witnessed  $k$ -distance.

**Lemma 4.** *If  $U$  is a point cloud with non-positive weights,  $\text{VR}_r(U) \subseteq \check{\text{C}}_{2r}(U)$ .*

*Proof.* Let  $u, v$  be two weighted points such that the balls  $\text{B}(u, (r^2 + w_u)^{1/2})$  and  $\text{B}(v, (r^2 + w_v)^{1/2})$  intersect. Let  $\ell$  denote the Euclidean distance between  $u$  and  $v$ . By hypothesis, we know that one of the two radii is at least  $\ell/2$ . Suppose  $w_u > w_v$ ; in this case,  $(r^2 + w_u)^{1/2} \geq \ell/2$ . Since the weights are non-positive, we also know that  $r \geq \ell/2$ . Using these two facts, we deduce

$$(2r)^2 + w_u = 3r^2 + (r^2 + w_u) \geq 3\ell^2/4 + \ell^2/4 = \ell^2.$$

This means that the point  $v$  belongs to the ball  $\text{B}(u, ((2r)^2 + w_u)^{1/2})$ .

Now, choose a simplex  $\sigma$  in the Vietoris-Rips complex  $\text{VR}_r(U)$ . Let  $v$  be its vertex with the smallest weight. By the previous paragraph, we know that  $v$  belongs to every ball  $\text{B}(u, (2r)^{1/2} + w_u)$ , for every  $u \in \sigma$ . Therefore, all these balls intersect, and, by definition,  $\sigma$  belongs to the Čech complex  $\check{\text{C}}_{2r}(U)$ .  $\square$

**Inference.** Suppose we are in the conditions of the Approximation Theorem. Additionally, we assume that the support  $K$  of the original measure  $\mu$  has a *weak feature size* larger than  $R$ . By definition of the weak feature size, the distance function  $d_K$  has no critical value in the interval  $(0, R)$ . Consequently, all the offsets  $K^r = d_K^{-1}[0, r]$  of  $K$  are homotopy equivalent for  $r \in (0, R)$ . Suppose again that we have drawn a set  $P$  of  $N$  points from a Wasserstein approximation  $\nu$  of  $\mu$ , such that  $W_2(\mu, \nu) \leq \sigma$ . From the Approximation Theorem, we have

$$\|d_{P,k}^w - d_K\|_\infty \leq e := 27m_0^{-1/2}\sigma + 12m_0^{1/\ell}\alpha_\mu^{-1/\ell}$$

with high probability as  $N$  goes to infinity. Then, the standard argument [10] shows that the Betti numbers of the compact set  $K$  can be inferred from the function  $d_{P,k}^w$ , which is defined only from the point sample  $P$ , as long as  $e$  is less than  $R/4$ . Indeed, denoting by  $K^r$  and  $U^r$  the  $r$ -sublevel sets of the functions  $d_K$  and  $d_{P,k}^w$ , the sequence of inclusions

$$K^0 \subseteq U^e \subseteq K^{2e} \subseteq U^{3e} \subseteq K^{4e}$$

holds with high probability. By assumption, the function  $d_K$  has no critical values in the range  $(0, 4e) \subseteq (0, R)$ . Therefore, the rank of the image on the homology induced by inclusion  $\text{H}(U^e) \rightarrow \text{H}(U^{3e})$  is equal to the Betti numbers of the set  $K$ . In the language of persistent homology [16], the persistent Betti numbers  $\beta^{(e, 3e)}$  of the function  $d_{P,k}^w$  are equal to the Betti numbers of the set  $K$ . Computationally, we can construct the sublevel sets  $U^e$  as either the Čech complex or the alpha shape.

Using the interleaving of Vietoris-Rips and Čech complexes, proved in Lemma 4, we can recover the Betti numbers from the Vietoris-Rips complex if  $e < R/9$  [7].

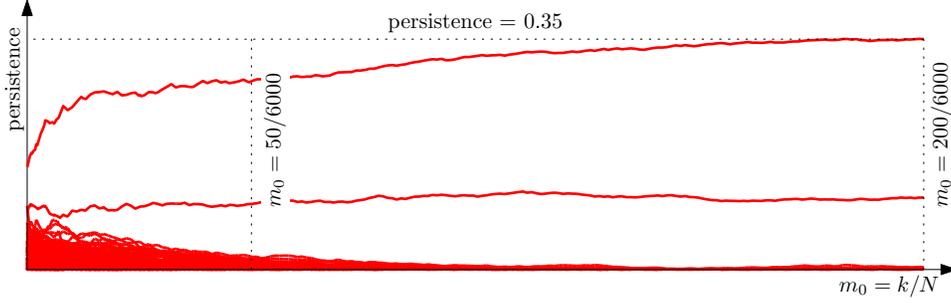


FIGURE 2. (PL-approximation of the) 1-dimensional persistence vineyard of the witnessed  $k$ -distance function. Topological features of the space, obscured by noise for low values of  $m_0$ , stand out as we increase the mass parameter.

From the following diagram of inclusions and homotopy equivalences

$$\begin{array}{ccccccc}
 \check{C}_e \subseteq \text{VR}_e \subseteq \check{C}_{2e} & \subseteq & \check{C}_{4e} \subseteq \text{VR}_{4e} \subseteq \check{C}_{8e} \\
 \wr & & \wr & & \wr & & \wr \\
 K_0 \subseteq U^e & \subseteq & U^{2e} \subseteq K^{3e} \subseteq U^{4e} & \subseteq & U^{8e} \subseteq K^{9e},
 \end{array}$$

it follows that the map on homology  $\text{H}(\text{VR}_e(U)) \rightarrow \text{H}(\text{VR}_{4e}(U))$  has the same rank as the homology of the space  $K$ .

**Choice of the mass parameter.** The language of persistent homology also suggests a strategy for choosing a mass parameter  $m_0$  for the distance to a measure — a question not addressed by the original paper [6]. For every mass parameter  $m_0$ , the  $p$ -dimensional *persistence diagram*  $\text{Pers}_p(d_{\mu, m_0})$  is a set of points  $\{(b_i(m_0), d_i(m_0))\}_i$  in the extended plane  $(\mathbb{R} \cup \{\infty\})^2$ . Each of these points represents a homology class of dimension  $p$  in the sublevel sets of  $d_{\mu, m_0}$ ;  $b_i(m_0)$  and  $d_i(m_0)$  are the values at which it is born and dies. Since the distance to measure  $d_{1_P, m_0}$  depends continuously on  $m_0$ , by the Stability Theorem [10] so do its persistence diagrams. Thus, one can use the vineyards algorithm [11] to track their evolution. Figure 2 illustrates such a construction for the point set in Figure 1 and the witnessed  $k$ -distance. It displays the evolution of the persistence  $(d_1(m_0) - b_1(m_0))$  of each of the 1-dimensional homology classes as  $m_0$  varies. This graph highlights the choices of the mass parameter that expose the two prominent classes (corresponding to the two loops of the figure-8).

#### ACKNOWLEDGEMENT

This work has been partly supported by a grant from the French ANR, ANR-09-BLAN-0331-01, NSF grants FODAVA 0808515, CCF 1011228, and NSF/NIH grant 0900700. The second author would also like to acknowledge the support of the Fields Institute during the redaction of this article.

#### REFERENCES

- [1] N. Amenta and M. Bern. Surface reconstruction by Voronoi filtering. *Discrete and Computational Geometry*, 22(4):481–504, 1999.
- [2] S. Arya and D. Mount. Computational geometry: proximity and location. *Handbook of Data Structures and Applications*, pages 63.1–63.22, 2005.

- [3] F. Aurenhammer. A New Duality Result Concerning Voronoi Diagrams. *Discrete and Computational Geometry*, 5(1):243–254, 1990.
- [4] F. Bolley, A. Guillin, and C. Villani. Quantitative Concentration Inequalities for Empirical Measures on Non-compact Spaces. *Probability Theory and Related Fields*, 137(3):541–593, 2007.
- [5] F. Chazal, D. Cohen-Steiner, and A. Lieutier. A sampling theory for compact sets in Euclidean space. *Discrete and Computational Geometry*, 41(3):461–479, 2009.
- [6] F. Chazal, D. Cohen-Steiner, and Q. Mérigot. Geometric inference for probability measures. Preprint (INRIA RR-6930 v2), 2010.
- [7] F. Chazal and S. Oudot. Towards persistence-based reconstruction in Euclidean spaces. *Proceedings of the ACM Symposium on Computational Geometry*, pages 232–241, 2008.
- [8] K. Clarkson. Nearest-neighbor searching and metric space dimensions. *Nearest-Neighbor Methods for Learning and Vision: Theory and Practice*, pages 15–59, 2006.
- [9] K. Clarkson and P. Shor. Applications of random sampling in computational geometry, II. *Discrete and Computational Geometry*, 4:387–421, 1989.
- [10] D. Cohen-Steiner, H. Edelsbrunner, and J. Harer. Stability of persistence diagrams. *Discrete and Computational Geometry*, 37(1):103–120, 2007.
- [11] D. Cohen-Steiner, H. Edelsbrunner, and D. Morozov. Vines and vineyards by updating persistence in linear time. In *Proceedings of the ACM Symposium on Computational Geometry*, pages 119–126, 2006.
- [12] S. Dasgupta. Learning mixtures of Gaussians. In *Proceedings of the IEEE Symposium on Foundations of Computer Science*, page 634, 1999.
- [13] T. Dey and S. Goswami. Provable surface reconstruction from noisy samples. *Computational Geometry*, 35(1-2):124–141, 2006.
- [14] D. Dubhashi, A. Panconesi, and C. U. Press. *Concentration of measure for the analysis of randomized algorithms*. Cambridge University Press, 2009.
- [15] H. Edelsbrunner. The union of balls and its dual shape. *Discrete and Computational Geometry*, 13:415–440, 1995.
- [16] H. Edelsbrunner and J. Harer. Persistent homology — a survey. *Surveys on Discrete and Computational Geometry. Twenty Years Later*, pages 257–282, 2008.
- [17] H. Edelsbrunner and J. Harer. *Computational Topology*. AMS, 2009.
- [18] P. Indyk. Nearest neighbors in high-dimensional spaces. *Handbook of Discrete and Computational Geometry*, pages 877–892, 2004.
- [19] B. Kloeckner. Approximation by finitely supported measures. Preprint (arXiv:1003.1035), 2010.
- [20] P. Niyogi, S. Smale, and S. Weinberger. A topological view of unsupervised learning from noisy data. Preprint, 2008.
- [21] P. Niyogi, S. Smale, and S. Weinberger. Finding the homology of submanifolds with high confidence from random samples. *Discrete and Computational Geometry*, 39(1):419–441, 2008.
- [22] C. Villani. *Topics in Optimal Transportation*. American Mathematical Society, 2003.

## APPENDIX A. PAIRWISE DISTANCES

A valuable property of the Vietoris-Rips complex is that the pairwise distances between the points suffice for its construction. We (re-)construct this property for the Vietoris-Rips complex on the witnessed barycenters. To do so, we need the weights of the barycenters as well as their pairwise distances in terms of the pairwise distances between the points of  $P$ .

**Intersection criteria.** Suppose we are given two weighted barycenters  $\bar{p} = (1/k) \sum_{i=1}^k p_i$  and  $\bar{q} = (1/k) \sum_{i=1}^k q_i$ . We start by finding the intersection point between the line  $(\bar{p}\bar{q})$  and the bisector of the Power cells of  $\bar{p}$  and  $\bar{q}$ . The point

$x_t = (1-t)\bar{p} + t\bar{q}$  belongs to this bisector if and only if:

$$\begin{aligned} \|x_t - \bar{p}\|^2 - w_{\bar{p}} &= \|x_t - \bar{q}\|^2 - w_{\bar{q}} \iff t^2\|\bar{p} - \bar{q}\|^2 - w_{\bar{p}} = (1-t)^2\|\bar{p} - \bar{q}\|^2 - w_{\bar{q}} \\ &\iff 2t = 1 + \frac{w_{\bar{p}} - w_{\bar{q}}}{\|\bar{p} - \bar{q}\|^2} \end{aligned}$$

The two balls  $B(\bar{p}, (r^2 + w_{\bar{p}})^{1/2})$  and  $B(\bar{q}, (r^2 + w_{\bar{q}})^{1/2})$  intersect if and only if the point  $x_t$  belongs to one of them, in which case it also belongs to the other. With the value of  $t$  that we found, this is equivalent to

$$\begin{aligned} \|x_t - \bar{p}\|^2 \leq r^2 + w_{\bar{p}} &\iff t^2\|\bar{p} - \bar{q}\|^2 - w_{\bar{p}} \leq r^2 \\ &\iff \frac{1}{4} \left( 1 + \frac{w_{\bar{p}} - w_{\bar{q}}}{\|\bar{p} - \bar{q}\|^2} \right)^2 \|\bar{p} - \bar{q}\|^2 - w_{\bar{p}} \leq r^2 \end{aligned}$$

Consequently, one can determine whether a segment  $\{\bar{p}, \bar{q}\}$  belongs to the Vietoris-Rips complex of the witnessed barycenters with parameter  $r$  by knowing only the weights of the barycenters and their pairwise distances. In the next two paragraphs, we show how to express these quantities in terms of the pairwise distances between the data points.

**Vertex weights.** For a barycenter  $\bar{p} = \frac{1}{k}(p_1 + \dots + p_k)$  of  $k$  distinct points of  $P$ ,

$$\begin{aligned} -w_{\bar{p}} &= \frac{1}{k} \sum_{i=1}^k \|\bar{p} - p_i\|^2 = \frac{1}{k} \sum_{i=1}^k \left\| \frac{1}{k} \sum_{j=1}^k p_j - p_i \right\|^2 = \frac{1}{k} \sum_{i=1}^k \left\| \frac{1}{k} \sum_{j=1}^k (p_j - p_i) \right\|^2 \\ &= \frac{1}{k^3} \sum_{i=1}^k \sum_{j=1}^k \sum_{l=1}^k \langle p_j - p_i | p_l - p_i \rangle \\ &= \frac{1}{k^3} \sum_{i=1}^k \sum_{j>i}^k \sum_{l>j}^k (\|p_i - p_j\|^2 + \|p_i - p_l\|^2 + \|p_j - p_l\|^2) \\ &= \frac{k-2}{k^3} \sum_{i=1}^k \sum_{j>i}^k \|p_i - p_j\|^2 \end{aligned}$$

The second to last equality is obtained by considering every triangle  $\triangle(p_i, p_j, p_l)$  and observing that

$$\begin{aligned} 0 &= (p_i - p_j + p_j - p_l + p_l - p_i)^2 \\ &= \|p_i - p_j\|^2 + \|p_j - p_l\|^2 + \|p_l - p_i\|^2 \\ &\quad + 2\langle p_i - p_j | p_j - p_l \rangle + 2\langle p_i - p_j | p_l - p_i \rangle + 2\langle p_j - p_l | p_l - p_i \rangle, \end{aligned}$$

and the last equality is obtained by observing that each edge appears in  $(k-2)$  triangles.

**Barycenter distances.** It remains to express the distance  $\|\bar{p} - \bar{q}\|^2$  between the barycenters in terms of the pairwise distances between the points  $\{p_1, \dots, p_k, q_1, \dots, q_k\}$ .

$$\|\bar{p} - \bar{q}\|^2 = \left\| \frac{1}{k} \left( \sum q_i - \sum p_i \right) \right\|^2 = \frac{1}{k^2} \left\| \sum (q_i - p_i) \right\|^2 = \frac{1}{k^2} \sum_{i=1}^k \sum_{j=1}^k \langle (q_i - p_i) | (q_j - p_j) \rangle$$

Rewriting

$$\langle (q_i - p_i) | (q_j - p_j) \rangle = \langle (q_i - p_i) | (q_j - p_i + p_i - p_j) \rangle = \langle (q_i - p_i) | (q_j - p_i) \rangle - \langle (q_i - p_i) | (p_j - p_i) \rangle$$

we get dot products between vectors with the same base point, which we express in terms of the areas of their respective triangles:

$$\langle (q_i - p_i) | (p_j - p_i) \rangle^2 = \|q_i - p_i\|^2 \|p_j - p_i\|^2 \cos^2 \theta = \|q_i - p_i\|^2 \|p_j - p_i\|^2 - 4S^2,$$

where  $S$  is the area of the triangle  $\triangle(p_i, q_i, p_j)$ . We compute it from the pairwise distances using Heron's formula,  $S^2 = s(s-a)(s-b)(s-c)$ , where  $s$  is the semiperimeter, and  $a, b, c$  are the lengths of the sides of the triangle.

*E-mail address:* `guibas@cs.stanford.edu`

DEPARTMENT OF COMPUTER SCIENCE, STANFORD UNIVERSITY

*E-mail address:* `dmitriy@mrzv.org`

LAWRENCE BERKELEY NATIONAL LABORATORY

*E-mail address:* `quentin.merigot@imag.fr`

LABORATOIRE JEAN KUNTZMANN, UNIVERSITÉ DE GRENOBLE AND CNRS