

# Topological Analysis of Scalar Fields with Outliers\*

Mickaël Buchet<sup>1</sup>, Frédéric Chazal<sup>1</sup>, Tamal K. Dey<sup>2</sup>, Fengtao Fan<sup>2</sup>,  
Steve Y. Oudot<sup>1</sup>, and Yusu Wang<sup>2</sup>

1 Inria Saclay Île-de-France, Palaiseau, France

mickael.buchet@m4x.org, frederic.chazal@inria.fr, steve.oudot@inria.fr

2 Department of Computer Science and Engineering, The Ohio State University,  
Columbus, OH 43210, USA

tamaldey@cse.ohio-state.edu, fan.171@osu.edu, yusu@cse.ohio-state.edu

---

## Abstract

Given a real-valued function  $f$  defined over a manifold  $M$  embedded in  $\mathbb{R}^d$ , we are interested in recovering structural information about  $f$  from the sole information of its values on a finite sample  $P$ . Existing methods provide approximation to the persistence diagram of  $f$  when geometric noise and functional noise are bounded. However, they fail in the presence of aberrant values, also called outliers, both in theory and practice.

We propose a new algorithm that deals with outliers. We handle aberrant functional values with a method inspired from the  $k$ -nearest neighbors regression and the local median filtering, while the geometric outliers are handled using the distance to a measure. Combined with topological results on nested filtrations, our algorithm performs robust topological analysis of scalar fields in a wider range of noise models than handled by current methods. We provide theoretical guarantees and experimental results on the quality of our approximation of the sampled scalar field.

**1998 ACM Subject Classification** I.3.5 Computational Geometry and Object Modeling

**Keywords and phrases** Persistent Homology, Topological Data Analysis, Scalar Field Analysis, Nested Rips Filtration, Distance to a Measure

**Digital Object Identifier** 10.4230/LIPIcs.SOCG.2015.827

## 1 Introduction

Consider a network of sensors measuring a quantity such as the temperature, the humidity, or the elevation. These sensors also compute their positions and communicate these data to others. However, they are not perfect and can make mistakes such as providing some aberrant values. Can we still recover topological structure from the measured quantity?

This is an instance of a scalar field analysis problem. Given a manifold  $M$  embedded in  $\mathbb{R}^d$  and a scalar field  $f : M \rightarrow \mathbb{R}$ , we want to extract topological information about  $f$ , knowing only its values on a finite set of points  $P$ . The critical points of a function, that is, peaks (local maxima), pits (local minima), and passes (saddle points) constitute important topological features of the function. In addition, the prominence of these features also contains valuable information, which the geographers use to distinguish between a summit and a local maximum in its shadow. Such information can be captured by the so-called *topological persistence*, which studies the *sub-level sets*  $f^{-1}((-\infty, \alpha])$  of a function  $f$  and the way their topology evolves as parameter  $\alpha$  increases. In the case of geography, we can use

---

\* See [1] for the full version of this paper.



the negated elevation as a function to study the topography. Peaks will appear depending on their altitude and will merge into other topological features at saddle points. This provides a *persistence diagram* describing the lifespan of features where the peaks with more prominence have longer lifespans.

When the domain  $M$  of the function  $f$  is triangulated, one classical way of computing this diagram is to linearly interpolate the function  $f$  on each simplex and then apply the standard persistence algorithm to this piecewise-linear function [16]. For cases where we only have pairwise distances between input points, one can build a family of simplicial complexes and infer the persistent homology of the input function  $f$  from them [6] (this construction will be detailed in Section 2).

Both of these approaches can provably approximate persistent homology when the input points admit a bounded noise, i.e., when the Hausdorff distance between  $P$  and  $M$  is bounded and the  $L_\infty$ -error on the observed value of  $f$  is also bounded. What happens if the noise is unbounded? A faulty sensor can provide completely wrong information or a bad position. Previous methods no longer work in this setting. Moreover, a sensor with a good functional value but a bad position can become an outlier in function value at its measured position (see Section 3.1 for an example). In this paper, we study the problem of analyzing scalar fields in the presence of unbounded noise both in the geometry and in the functional values. To the best of our knowledge, there is no other method to handle such combined unbounded geometric and functional noise with theoretical guarantees.

### Contributions

We consider a general sampling condition. Intuitively, a sample  $(P, \tilde{f})$  of a function  $f : M \rightarrow \mathbb{R}$  respects our condition if: (i) the domain  $M$  is sampled densely and there is no cluster of noisy samples outside  $M$  (roughly speaking, no area outside  $M$  has a higher sampling density than on  $M$ ), and (ii) for any point of  $P$ , at least half of its  $k$  nearest neighbors have a functional value with an error less than a threshold  $s$ . This condition allows functional outliers that may have a value arbitrarily far away from the true one. It encompasses the previous bounded sampling conditions as well as other sampling conditions such as bounded Wasserstein distance for geometry, or generative models like an additive Gaussian noise. Connection to some of these classical sampling conditions can be found in the full version of the paper [1].

We show how to approximate the persistence diagram of  $f$  knowing only its observed value  $\tilde{f}$  on the set  $P$ . We achieve this goal through three main steps:

1. Using the observations  $\tilde{f}$ , we provide a new estimator  $\hat{f}$  to approximate  $f$ . This estimator is inspired by the  $k$ -nearest neighbours regression technique but differs from it in an essential way.
2. We filter geometric outliers using a distance to a measure function.
3. We combine both techniques in a unified framework to estimate the persistence diagram of  $f$ .

The two sources of noise, geometric and functional, are not independent. The interdependency is first identified by assuming appropriate sampling conditions, and then untangled by separate steps in our algorithm.

### Related work

A framework for scalar field topology inference with theoretical guarantees has been previously proposed in [6]. However, it is limited to a bounded noise assumption, which we aim to relax.

For handling the functional noise only, the traditional non-parametric regression mostly uses kernel-based or  $k$ -NN estimators. The  $k$ -NN methods are more versatile [13]. Nevertheless, the kernel-based estimators are preferred when there is structure in the data. However, the functional outliers destroy the structure on which kernel-based estimators rely. These functional outliers can arise as a result of geometric outliers (see Section 3.1). Thus, in a way, it is essential to be able to handle functional outliers when the input has geometric noise. Functional outliers can also introduce a bias that hampers the robustness of a  $k$ -NN regression. For example, if all outliers' values are greater than the actual value, a  $k$ -NN regression will shift towards a larger value. Our approach leverages the  $k$ -NN regression idea while trying to avoid the sensitivity to this bias.

Various methods for geometric denoising have also been proposed in the literature. If the generative model for noise is known a priori, one can use de-convolution to remove noise. Some methods have been specifically adapted to use topological information for such denoising [14]. In our case where the generative model is unknown, we use a filtering by the value of the distance to a measure, which has been successfully applied to infer the topology of a domain under unbounded noise [4].

## 2 Preliminaries for Scalar Field Analysis

In [6], Chazal et al. presented an algorithm to analyze the scalar field topology using persistent homology which can handle bounded Hausdorff noise both in geometry and in observed function values. Our approach follows the same high level framework. Hence in this section, we introduce necessary preliminaries along with some of the results from [6].

### Riemannian manifold and its sampling.

Consider a compact Riemannian manifold  $M$ . Let  $d_M$  denote the geodesic metric on  $M$ . Consider the open Riemannian ball  $B_M(x, r) := \{y \in M \mid d_M(x, y) < r\}$  centered at  $x \in M$ .  $B_M(x, r)$  is *strongly convex* if for any pair  $(y, y')$  in the closure of  $B_M(x, r)$ , there exists a unique minimizing geodesic between  $y$  and  $y'$  whose interior is contained in  $B_M(x, r)$ . Given any  $x \in M$ , let  $\varrho(x)$  denote the supremum of the value of  $r$  such that  $B_M(x, r)$  is strongly convex. As  $M$  is compact, the infimum of all  $\varrho(x)$  is positive and we denote it by  $\varrho(M)$ , which is called the *strong convexity radius* of  $M$ .

A point set  $P \subseteq M$  is a *geodesic  $\varepsilon$ -sample* of  $M$  if for every point  $x$  of  $M$ , the distance from  $x$  to  $P$  is less than  $\varepsilon$  in the metric  $d_M$ . Given a  $c$ -Lipschitz scalar function  $f : M \rightarrow \mathbb{R}$ , we aim to study the persistent homology of  $f$ . However, the scalar field  $f : M \rightarrow \mathbb{R}$  is only approximated by a discrete set of sample points  $P$  and a function  $\tilde{f} : P \rightarrow \mathbb{R}$ . The goal of this paper is to retrieve the topological structure of  $f$  from  $\tilde{f}$  when some forms of noise are present both in the positions of  $P$  and in the function values of  $\tilde{f}$ .

### Persistent homology.

As in [6], we infer the persistent homology of  $f$  using well-chosen *persistence modules*. A *filtration*  $\{F_\alpha\}_{\alpha \in \mathbb{R}}$  is a family of sets  $F_\alpha$  totally ordered by inclusions  $F_\alpha \subseteq F_\beta$ . Following [3], a persistence module is a family of vector spaces  $\{\Phi_\alpha\}_{\alpha \in \mathbb{R}}$  with a family of homomorphisms  $\phi_\alpha^\beta : \Phi_\alpha \rightarrow \Phi_\beta$  such that for all  $\alpha \leq \beta \leq \gamma$ ,  $\phi_\alpha^\gamma = \phi_\beta^\gamma \circ \phi_\alpha^\beta$ . Given a filtration  $\mathcal{F} = \{F_\alpha\}_{\alpha \in \mathbb{R}}$  and  $\alpha \leq \beta$ , the canonical inclusion  $F_\alpha \hookrightarrow F_\beta$  induces a homomorphism at the homology level  $H_*(F_\alpha) \rightarrow H_*(F_\beta)$ . These homomorphisms and the homology groups of  $F_\alpha$  form the so-called *persistence module* of  $\mathcal{F}$ .

The persistence module of the filtration  $\mathcal{F} = \{F_\alpha\}_{\alpha \in \mathbb{R}}$  is said to be *q-tame* when all the homomorphisms  $H_*(F_\alpha) \rightarrow H_*(F_\beta)$  have finite rank [5]. Its algebraic structure can then be described by the *persistence diagram*  $\text{Dgm}(\mathcal{F})$ , which is a multiset of points in  $\mathbb{R}^2$  describing the lifespan of the homological features in the filtration  $\mathcal{F}$ . For technical reasons,  $\text{Dgm}(\mathcal{F})$  also contains every point of the diagonal  $y = x$  with countably infinite multiplicity. See [10] for a more formal discussion of the persistence diagrams.

Persistence diagrams can be compared using the *bottleneck distance*  $d_B$  [8]. Given two multisets with the same cardinality, possibly infinite,  $D$  and  $E$  in  $\mathbb{R}^2$ , we consider the set  $\mathcal{B}$  of all bijections between  $D$  and  $E$ . The bottleneck distance (under  $L_\infty$ -norm) is then defined as:

$$d_B(D, E) = \inf_{b \in \mathcal{B}} \sup_{x \in D} \|x - b(x)\|_\infty. \tag{1}$$

Two filtrations  $\{U_\alpha\}$  and  $\{V_\alpha\}$  are said to be  $\varepsilon$ -interleaved if, for any  $\alpha$ , we have  $U_\alpha \subset V_{\alpha+\varepsilon} \subset U_{\alpha+2\varepsilon}$ . Recent work in [3, 5] shows that two interleaved filtrations induce close persistence diagrams in the bottleneck distance.

► **Theorem 2.1.** *Let  $U$  and  $V$  be two  $q$ -tame and  $\varepsilon$ -interleaved filtrations. Then the persistence diagrams of these filtrations verify  $d_B(\text{Dgm}(U), \text{Dgm}(V)) \leq \varepsilon$ .*

**Nested filtrations**

The scalar field topology of  $f : M \rightarrow \mathbb{R}$  is studied via the topological structure of the sub-level sets filtration of  $f$ . More precisely, the sub-level sets of  $f$  are defined as  $F_\alpha = f^{-1}((-\infty, \alpha])$  for any  $\alpha \in \mathbb{R}$ . The collection of sub-level sets forms a filtration  $\mathcal{F} = \{F_\alpha\}_{\alpha \in \mathbb{R}}$  connected by natural inclusions  $F_\alpha \subseteq F_\beta$  for any  $\alpha \leq \beta$ . Our goal is to approximate the persistence diagram  $\text{Dgm}(\mathcal{F})$  from the observed scalar field  $\tilde{f} : P \rightarrow \mathbb{R}$ . We now describe the results of [6] for approximating  $\text{Dgm}(\mathcal{F})$  when  $P$  is a geodesic  $\varepsilon$ -sample of  $M$ . These results will later be useful for our approach.

To simulate the sub-level sets filtration  $\{F_\alpha\}$  of  $f$ , we introduce  $P_\alpha = \tilde{f}^{-1}((-\infty, \alpha]) \subseteq P$  for any  $\alpha \in \mathbb{R}$ . The points in  $P_\alpha$  intuitively sample the sub-level set  $F_\alpha$ . To estimate the topology of  $F_\alpha$  from these discrete samples  $P_\alpha$ , we consider the  $\delta$ -offset  $P^\delta$  of the point set  $P$ , i.e., we grow geodesic balls of radius  $\delta$  around the points of  $P$ . This gives us a union of balls that serves as a proxy for  $f^{-1}((-\infty, \alpha])$ . The nerve of this collection of balls, also known as the *Čech complex*,  $C_\delta(P)$ , has many interesting properties but is difficult to compute in high dimensions. We consider an alternate complex called the *Vietoris-Rips complex*  $R_\delta(P)$  that is easier to compute. It is defined as the maximal simplicial complex with the same 1-skeleton as the Čech complex. The Čech and Rips complexes are related in any metric space:  $\forall \delta > 0, C_\delta(P) \subset R_\delta(P) \subset C_{2\delta}(P)$ .

Even though a single Vietoris-Rips complex may not capture the homology of the manifold  $M$ , a pair of nested complexes can recover it using the inclusions  $R_\delta(P_\alpha) \hookrightarrow R_{2\delta}(P_\alpha)$  [7]. Specifically, for a fixed  $\delta > 0$ , consider the following commutative diagram induced by inclusions, for  $\alpha \leq \beta$ :

$$\begin{array}{ccc} H_*(R_{2\delta}(P_\alpha)) & \xrightarrow{\phi_\alpha^\beta} & H_*(R_{2\delta}(P_\beta)) \\ i_\alpha \uparrow & & \uparrow i_\beta \\ H_*(R_\delta(P_\alpha)) & \longrightarrow & H_*(R_\delta(P_\beta)) \end{array}$$

As the diagram commutes for all  $\alpha \leq \beta$ ,  $\{Im(i_\alpha), \phi_\alpha^\beta|_{Im(i_\alpha)}\}$  defines a persistence module. We call it the persistent homology module of the filtration of nested pairs  $\{R_\delta(P_\alpha) \hookrightarrow$

$R_{2\delta}(P_\alpha)\}_{\alpha \in \mathbb{R}}$ . This construction can also be done for any filtration of nested pairs. Using this construction, one of the main results of [6] is:

► **Theorem 2.2** (Theorems 2 and 6 of [6]). *Let  $M$  be a compact Riemannian manifold and let  $f : M \rightarrow \mathbb{R}$  be a  $c$ -Lipschitz function. Let  $P$  be a geodesic  $\varepsilon$ -sample of  $M$ . If  $\varepsilon < \frac{1}{4}\varrho(M)$ , then for any  $\delta \in [2\varepsilon, \frac{1}{2}\varrho(M))$ , the persistent homology modules of  $f$  and of the filtration of nested pairs  $\{R_\delta(P_\alpha) \hookrightarrow R_{2\delta}(P_\alpha)\}$  are  $2c\delta$ -interleaved. Therefore, the bottleneck distance between their persistence diagrams is at most  $2c\delta$ .*

*Furthermore, the  $k$ -dimensional persistence diagram for the filtrations of nested pairs  $\{R_\delta(P_\alpha) \hookrightarrow R_{2\delta}(P_\alpha)\}$  can be computed in  $O(|P|kN + N \log N + N^3)$  time, where  $N$  is the number of simplices of  $\{R_{2\delta}(P_\infty)\}$ , and  $|P|$  denotes the cardinality of the sample set  $P$ .*

It has been observed that, in practice, the persistence algorithm often has a running time linear in the number of simplices, which reduces the above complexity to  $O(|P| + N \log N)$  in a practical setting.

We say that  $\tilde{f}$  has a precision of  $\xi$  over  $P$  if  $|\tilde{f}(p) - f(p)| \leq \xi$  for any  $p \in P$ . We then have the following result for the case when we only have this functional noise:

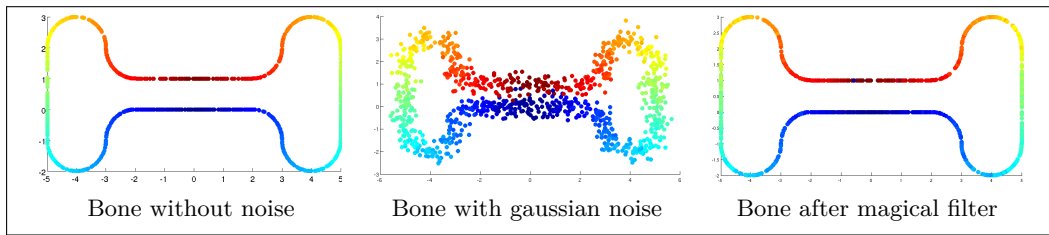
► **Theorem 2.3** (Theorem 3 of [6]). *Let  $M$  be a compact Riemannian manifold and let  $f : M \rightarrow \mathbb{R}$  be a  $c$ -Lipschitz function. Let  $P$  be a geodesic  $\varepsilon$ -sample of  $M$  such that the values of  $f$  on  $P$  are known with precision  $\xi$ . If  $\varepsilon < \frac{1}{4}\varrho(M)$ , then for any  $\delta \in [2\varepsilon, \frac{1}{2}\varrho(M))$ , the persistent homology modules of  $f$  and of the filtration of nested pairs  $\{R_\delta(P_\alpha) \hookrightarrow R_{2\delta}(P_\alpha)\}$  are  $(2c\delta + \xi)$ -interleaved. Therefore, the bottleneck distance between their persistence diagrams is at most  $2c\delta + \xi$ .*

Geometric noise was considered in the form of bounded noise in the estimate of the geodesic distances between points in  $P$ . It translated into a relation between the measured pairwise distances and the real ones. With only geometric noise, one has the following stability result. It was stated in this form in the conference version of the paper.

► **Theorem 2.4** (Theorem 4 of [6]). *Let  $M, f$  be defined as previously and  $P$  be an  $\varepsilon$ -sample of  $M$  in its Riemannian metric. Assume that, for a parameter  $\delta > 0$ , the Rips complexes  $R_\delta(\cdot)$  are defined with respect to a metric  $\tilde{d}(\cdot, \cdot)$  which satisfies  $\forall x, y \in P, \frac{d_M(x, y)}{\lambda} \leq \tilde{d}(x, y) \leq \nu + \mu \frac{d_M(x, y)}{\lambda}$ , where  $\lambda \geq 1$  is a scaling factor,  $\mu \geq 1$  is a relative error and  $\nu \geq 0$  an additive error. Then, for any  $\delta \geq \nu + 2\mu\frac{\varepsilon}{\lambda}$  and any  $\delta' \in [\nu + 2\mu\delta, \frac{1}{\lambda}\varrho(M)]$ , the persistent homology modules of  $f$  and of the filtration of nested pairs  $\{R_\delta(P_\alpha) \hookrightarrow R_{\delta'}(P_\alpha)\}$  are  $c\lambda\delta'$ -interleaved. Therefore, the bottleneck distance between their persistence diagrams is at most  $c\lambda\delta'$ .*

### 3 Functional Noise

In this section, we focus on the case where we have only functional noise in the observed function  $\tilde{f}$ . Suppose we have a scalar function  $f$  defined on a Riemannian manifold  $M$  embedded in  $\mathbb{R}^d$ . Note that the results of section 3 hold if  $\mathbb{R}^d$  is replaced by a metric space  $\mathbb{X}$ . We are given a geodesic  $\varepsilon$ -sample  $P \subset M$ , and a noisy observed function  $\tilde{f} : P \rightarrow \mathbb{R}$ . Our goal is to approximate the persistence diagram  $\text{Dgm}(\mathcal{F})$  of the sub-level set filtration  $\mathcal{F} = \{F_\alpha = f^{-1}((-\infty, \alpha])\}_\alpha$  from  $\tilde{f}$ . We assume that  $f$  is  $c$ -Lipschitz with respect to the intrinsic metric of the Riemannian manifold  $M$ . Note that this does not imply a Lipschitz condition on  $\tilde{f}$ .



■ **Figure 1** Bone example after applying Gaussian perturbation and magical filter

### 3.1 Functional sampling condition

Previous work on functional noise focused on bounded noise (e.g, [6]) or noise with zero-mean (e.g, [15]). However, there are many practical scenarios where the observed function  $\tilde{f}$  may contain these previously considered types of noise combined with *aberrant function values* in  $\tilde{f}$ . Hence, we propose below a more general sampling condition that allows such combinations.

#### Motivating examples

First, we provide some motivating examples for the need of handling *aberrant* function values in  $\tilde{f}$ , where  $\tilde{f}(p)$  at some sample point  $p$  can be totally unrelated to the true value  $f(p)$ . Consider a sensor network, where each node returns some measures. Such measurements can be imprecise, and in addition to that, a sensor may experience failure and return a completely wrong measure that has no relation with the true value of  $f$ . Similarly, an image could be corrupted with impulse noise where there are random pixels with aberrant function values, such as random white or black dots.

More interestingly, outliers in function values can naturally appear as a result of (extrinsic) geometric noise present in the discrete samples. For example, imagine that we have a process that can measure the function value  $f : M \rightarrow \mathbb{R}$  with *no error*. However, the geometric location  $\tilde{p}$  of a point  $p \in M$  can be wrong. In particular,  $\tilde{p}$  can be close to other parts of the manifold, thereby although  $\tilde{p}$  has the correct function value  $f(p)$ , it becomes a functional outlier among its neighbors (due to the wrong location of  $\tilde{p}$ ). See Figure 1 for an illustration. The function defined on this bone-like curve is the geodesic distance to a base point. The two sides of the narrow neck have very different function values. Now, suppose that the points are sampled uniformly on  $M$  and their position is then perturbed by an additive Gaussian noise. Then, points from one side of this neck can be sent closer to the other side, causing aberrant values in the observed function.

In fact, even if we assume that we have a “magic filter” that can project each sample back to the closest point on the underlying manifold  $M$ , the result is a new set of samples where all points are on the manifold and thus can be seen as having **no** geometric noise; however, this point set now contains functional noise which is actually caused by the original geometric noise. Note that such a magic filter is the goal of many geometric denoising methods. A perfect algorithm in this sense cannot remove or may even cause more aberrant functional noise. This motivates the need for handling functional outliers (in addition to traditional functional noise) as well as processing noise that combines geometric and functional noise together and that does not necessarily have zero-mean.

Another case where our approach is useful concerns with missing data. Assuming that some of the functional values are missing, we can replace them by anything and act as if they were outliers. Without modifying the algorithm, we obtain a way to handle the local loss of information.

**Functional sampling condition**

To allow both aberrant and more traditional functional noise, we introduce the following sampling condition. Let  $P \subset M$  be a geodesic  $\varepsilon$ -sample of the underlying manifold  $M$ . Intuitively, our sampling condition requires that for every point  $p \in P$ , locally there is a sufficient number of sample points with reasonably good function values. Specifically, we fix two parameters  $k$  and  $k'$  with the condition that  $k \geq k' > \frac{1}{2}k$ . Let  $NN_P^k(p)$  denote the set of the  $k$ -nearest neighbors of  $p$  in  $P$  in the *extrinsic metric*. We say that a discrete scalar field  $\tilde{f} : P \rightarrow \mathbb{R}$  is a  $(k, k', \Delta)$ -functional-sample of  $f : M \rightarrow \mathbb{R}$  if the following holds:

$$\forall p \in P, \left| \left\{ q \in NN_P^k(p) \mid |\tilde{f}(q) - f(p)| \leq \Delta \right\} \right| \geq k' \tag{2}$$

Intuitively, this sampling condition allows up to  $k - k'$  samples around a point  $p$  to be outliers (whose function values deviates from  $f(p)$  by at least  $\Delta$ ). In the full version [1], we consider two standard functional sampling conditions used in the statistical learning community and look at what they correspond to in our setting.

**3.2 Functional Denoising**

Given a scalar field  $\tilde{f} : P \rightarrow \mathbb{R}$  which is a  $(k, k', \Delta)$ -functional-sample of  $f : M \rightarrow \mathbb{R}$ , we now aim to compute a denoised function  $\hat{f} : P \rightarrow \mathbb{R}$  from the observed function  $\tilde{f}$ , and we will later use  $\hat{f}$  to infer the topology of  $f : M \rightarrow \mathbb{R}$ . Below we describe two ways to denoise the noisy observation  $\tilde{f}$ : one of which is well-known, and the other one is new. As we will see later, these two treatments lead to similar theoretical guarantees in terms of topology inference. However, they have different characteristics in practice, which are discussed in the full version [1].

***k*-median denoising**

In the *k*-median treatment, we simply perform the following: given any point  $p \in P$ , we set  $\hat{f}(p)$  to be the median value of the set of  $\tilde{f}$  values for the  $k$ -nearest neighbors  $NN_P^k(p) \subseteq P$  of  $p$ . We call  $\hat{f}$  the *k*-median denoising of  $\tilde{f}$ . The following observation is straightforward:

► **Observation 1.** *If  $\tilde{f} : P \rightarrow \mathbb{R}$  is a  $(k, k', \Delta)$ -functional-sample of  $f : M \rightarrow \mathbb{R}$  with  $k' \geq k/2$ , then we have  $|\hat{f}(p) - f(p)| \leq \Delta$  for any  $p \in P$ , where  $\hat{f}$  is the *k*-median denoising of  $\tilde{f}$ .*

**Disparity-based denoising**

In the *k*-median treatment, we choose a single value from the  $k$ -nearest neighbors of a sample point  $p$  and set it to be the denoised value  $\hat{f}(p)$ . This value, while within  $\Delta$  distance to the true value  $f(p)$  for  $k' \geq k/2$ , tends to have greater variability among neighboring sample points. Intuitively, taking the average (such as *k*-means) makes the function  $\hat{f}(p)$  smoother, but it is sensitive to outliers. We combine these ideas together, and use the following concept of disparity to help us identify a subset of points from the *k*-nearest neighbors of a sample point  $p$  to estimate  $\hat{f}(p)$ .

Given a set  $Y = \{x_1, \dots, x_l\}$  of  $l$  sample points from  $P$ , we define its disparity w.r.t.  $\tilde{f}$  as:

$$\phi(Y) = \frac{1}{l} \sum_{i=1}^l (\tilde{f}(x_i) - \mu(Y))^2, \quad \text{where } \mu(Y) = \frac{1}{l} \sum_{i=1}^l \tilde{f}(x_i).$$

$\mu(Y)$  and  $\phi(Y)$  are respectively the average and the variance of the observed function values for points from  $Y$ . Intuitively,  $\phi(Y)$  measures how tight the function values ( $\tilde{f}(x_i)$ ) are clustered. Now, given a point  $p \in P$ , we define

$$\widehat{Y}_p = \operatorname{argmin}_{Y \subseteq \operatorname{NN}_P^k(p), |Y|=k'} \phi(Y), \quad \text{and} \quad \widehat{z}_p = \mu(\widehat{Y}_p).$$

That is,  $\widehat{Y}_p$  is the subset of  $k'$  points from the  $k$ -nearest neighbors of  $p$  that has the smallest disparity and  $\widehat{z}_p$  is its mass center. It turns out that  $\widehat{Y}_p$  and  $\widehat{z}_p$  can be computed by the following *sliding-window* procedure: (i) Sort  $\operatorname{NN}_P^k(p) = \{x_1, \dots, x_k\}$  according to  $\tilde{f}(x_i)$ . (ii) For every  $k'$  consecutive points  $Y_i = \{x_i, \dots, x_{i+k'-1}\}$  with  $i \in [1, k - k' + 1]$ , compute its disparity  $\phi(Y_i)$ . (iii) Set  $\widehat{Y}_p = \operatorname{argmin}_{Y_i, i \in [1, k - k' + 1]} \phi(Y_i)$ , and return  $\mu(\widehat{Y}_p)$  as  $\widehat{z}_p$ . In the *disparity-based denoising* approach, we simply set  $\widehat{f}(p) := \widehat{z}_p$  as computed above. The approximation guarantee of  $\widehat{f}$  for the function  $f$  is given by the following Lemma.

► **Lemma 3.1.** *If  $\tilde{f} : P \rightarrow \mathbb{R}$  is a  $(k, k', \Delta)$ -functional-sample of  $f : M \rightarrow \mathbb{R}$  with  $k' \geq \frac{k}{2}$ , then we have  $|\widehat{f}(p) - f(p)| \leq \left(1 + 2\sqrt{\frac{k-k'}{2k'-k}}\right) \Delta$  for every  $p \in P$ , where  $\widehat{f}$  is the disparity-based denoising of  $\tilde{f}$ . In particular, if  $k' \geq \frac{2}{3}k$ , then  $|\widehat{f}(p) - f(p)| \leq 3\Delta$  for every  $p \in P$ .*

**Proof.** Let  $Y_\Delta = \{x \in \operatorname{NN}_P^k(p) : |\tilde{f}(x) - f(p)| \leq \Delta\}$  be the set of points in  $\operatorname{NN}_P^k(p)$  whose observed function values are within distance  $\Delta$  from  $f(p)$ . Since  $\tilde{f}$  is a  $(k, k', \Delta)$ -functional-sample of  $f$ , it is clear that  $|Y_\Delta| \geq k'$ . Let  $Y'_\Delta \subset Y_\Delta$  be a subset with  $k'$  elements,  $Y'_\Delta = \{x'_i\}_{i=1}^{k'}$ . By the definitions of  $Y_\Delta$  and  $Y'_\Delta$ , one can immediately check that  $|\tilde{f}(x'_i) - \mu(Y'_\Delta)| \leq 2\Delta$  where  $\mu(Y'_\Delta) = \frac{1}{k'} \sum_{i=1}^{k'} \tilde{f}(x'_i)$ . This inequality then gives an upper bound of the disparity  $\phi(Y'_\Delta)$ ,

$$\begin{aligned} \phi(Y'_\Delta) &= \frac{1}{k'} \sum_{i=1}^{k'} (\tilde{f}(x'_i) - \mu(Y'_\Delta))^2 \\ &\leq \frac{1}{k'} \sum_{i=1}^{k'} (2\Delta)^2 \\ &= 4\Delta^2 \end{aligned}$$

Recall from the sliding window procedure that  $\widehat{Y}_p = \operatorname{argmin}_{Y_i, i \in [1, k - k' + 1]} \phi(Y_i)$  and  $\widehat{z}_p = \mu(\widehat{Y}_p)$ . Denote  $A_1 = \widehat{Y}_p \cap Y_\Delta$  and  $A_2 = \widehat{Y}_p \setminus A_1$ . Since  $\tilde{f}$  is a  $(k, k', \Delta)$ -functional-sample of  $f$ , the size of  $A_2$  is at most  $k - k'$  and  $|A_1| \geq 2k' - k$ . If  $|\widehat{z}_p - f(p)| \leq \Delta$ , nothing needs to be proved. Without loss of generality, one can assume that  $f(p) + \Delta \leq \widehat{z}_p$ . Denote  $\delta = \widehat{z}_p - (f(p) + \Delta)$ . The disparity of  $\phi(\widehat{Y}_p)$  can then be estimated.

$$\begin{aligned} \phi(\widehat{Y}_p) &= \frac{1}{k'} \left( \sum_{x \in A_1} (\tilde{f}(x) - \widehat{z}_p)^2 + \sum_{x \in A_2} (\tilde{f}(x) - \widehat{z}_p)^2 \right) \\ &\geq \frac{1}{k'} \left( |A_1| \delta^2 + \sum_{x \in A_2} (\tilde{f}(x) - \widehat{z}_p)^2 \right) \\ &\geq \frac{1}{k'} \left( |A_1| \delta^2 + \frac{1}{|A_2|} \left( \sum_{x \in A_2} \tilde{f}(x) - |A_2| \widehat{z}_p \right)^2 \right) \\ &= \frac{1}{k'} \left( |A_1| \delta^2 + \frac{1}{|A_2|} \left( \sum_{x \in A_1} \tilde{f}(x) - |A_1| \widehat{z}_p \right)^2 \right) \\ &\geq \frac{1}{k'} \left( |A_1| \delta^2 + \frac{1}{|A_2|} (|A_1| \delta)^2 \right) \\ &= \frac{1}{k'} \delta^2 \left( \frac{|A_1|}{|A_2|} (|A_1| + |A_2|) \right) \\ &\geq \frac{1}{k'} \delta^2 \left( \frac{k' |A_1|}{|A_2|} \right) \\ &\geq \frac{2k' - k}{k - k'} \delta^2 \end{aligned}$$

where the third line uses the inequality  $\sum_{i=1}^n a_i^2 \geq \frac{1}{n} (\sum_{i=1}^n a_i)^2$ , and the fourth line uses the fact that  $(|A_1| + |A_2|) \widehat{z}_p = \sum_{x \in \widehat{Y}_p} \tilde{f}(x)$ . Since  $\widehat{Y}_p = \operatorname{argmin}_{Y_i, i \in [1, k - k' + 1]} \phi(Y_i)$ , it holds that



$\phi(\widehat{Y}_p) \leq \phi(Y'_\Delta)$ . Therefore,

$$\frac{2k' - k}{k - k'} \delta^2 \leq 4\Delta^2.$$

It then follows that  $\delta \leq 2\sqrt{\frac{k-k'}{2k'-k}}\Delta$  and  $|\widehat{f}(p) - f(p)| \leq \left(1 + 2\sqrt{\frac{k-k'}{2k'-k}}\right)\Delta$  since  $\widehat{z}_p = \widehat{f}(p)$ . If  $k' \geq \frac{2}{3}k$ , then  $1 + 2\sqrt{\frac{k-k'}{2k'-k}} \leq 1 + 2 = 3$ , meaning that  $|\widehat{f}(p) - f(p)| \leq 3\Delta$  in this case. ◀

► **Corollary 3.2.** *Given a  $(k, k', \Delta)$ -functional-sample of  $f : M \rightarrow \mathbb{R}$  with  $k' \geq k/2$ , we can compute a new function  $\widehat{f} : P \rightarrow \mathbb{R}$  such that  $|\widehat{f}(p) - f(p)| \leq \xi\Delta$  for any  $p \in P$ , where  $\xi = 1$  under  $k$ -median denoising, and  $\xi = \left(1 + 2\sqrt{\frac{k-k'}{2k'-k}}\right)$  under the disparity-based denoising.*

Hence after the  $k$ -median denoising or the disparity-based denoising, we obtain a new function  $\widehat{f}$  whose value at each sample point is within  $\xi\Delta$  precision to the true function value. We can now apply the scalar field topology inference framework from [6] (as introduced in Section 2) using  $\widehat{f}$  as input. In particular, set  $L_\alpha = \{p \in P \mid \widehat{f}(p) \leq \alpha\}$ , and let  $R_\delta(X)$  denote the Rips complex over points in  $X$  with parameter  $\delta$ . We approximate the persistence diagram induced by the sub-level sets filtration of  $f : M \rightarrow \mathbb{R}$  from the filtrations of nested pairs  $\{R_\delta(L_\alpha) \hookrightarrow R_{2\delta}(L_\alpha)\}_\alpha$ . It follows from Theorem 2.3 that:

► **Theorem 3.3.** *Let  $M$  be a compact Riemannian manifold and let  $f : M \rightarrow \mathbb{R}$  be a  $c$ -Lipschitz function. Let  $P$  be a geodesic  $\varepsilon$ -sample of  $M$ , and  $\widehat{f} : P \rightarrow \mathbb{R}$  a  $(k, k', \Delta)$ -functional-sample of  $f$ . Set  $\xi = 1$  if  $P_\alpha$  is obtained via  $k$ -median denoising, and  $\xi = \left(1 + 2\sqrt{\frac{k-k'}{2k'-k}}\right)$  if  $P_\alpha$  is obtained via disparity-based denoising. If  $\varepsilon < \frac{1}{4}\varrho(M)$ , then for any  $\delta \in [2\varepsilon, \frac{1}{2}\varrho(M))$ , the persistent homology modules of  $f$  and the filtration of nested pairs  $\{R_\delta(P_\alpha) \hookrightarrow R_{2\delta}(P_\alpha)\}$  are  $(2c\delta + \xi\Delta)$ -interleaved. Therefore, the bottleneck distance between their persistence diagrams is at most  $2c\delta + \xi\Delta$ .*

The above theoretical results are similar for  $k$ -median and disparity-based methods with a slight advantage for the  $k$ -median. However, interesting experimental results can be obtained when the Lipschitz condition on the function is removed, for example with images, where the disparity based method appears to be more resilient to large amounts of noise than the  $k$ -median denoising method. Illustrating examples can be found in the full version [1].

## 4 Geometric noise

In the previous section, we assumed that we have no geometric noise in the input. In this section, we deal with the case where there is only geometric noise in the input, but no functional noise of any kind. Specifically, for any point  $p \in P$ , we assume that the observed value  $\tilde{f}(p)$  is equal to the true function value  $f(\pi(p))$  where  $\pi(p)$  is the nearest point projection of  $p$  to the manifold. If  $p$  is on the medial axis of  $M$ , the projection  $\pi$  is arbitrary to one of the nearest points. As we have alluded before, general geometric noise implicitly introduces functional noise because the point  $p$  may have become a functional aberration of its orthogonal projection  $\pi(p) \in M$ . This error will be ultimately dealt with in Section 5 when we combine the results on purely functional noise from the previous section with the results on purely geometric noise in this section.

### 4.1 Sampling condition

#### Distance to a measure

The distance to a measure is a tool introduced to deal with geometrically noisy datasets, which are modelled as probability measures [4]. Given a probability measure  $\mu$  on a metric space  $\mathbb{X}$ , we define the *pseudo-distance*  $\delta_m(x)$  for any point  $x \in \mathbb{R}^d$  and a mass parameter  $m \in (0, 1]$  as  $\delta_m(x) = \inf\{r \in \mathbb{R} \mid \mu(B(x, r)) \geq m\}$ . The distance to a measure is then defined by averaging this quantity:

$$d_{\mu,m}(x) = \sqrt{\frac{1}{m} \int_0^m \delta_l(x)^2 dl}.$$

The *Wasserstein distance* is a standard tool to compare two measures. Given two probability measures  $\mu$  and  $\nu$  on a metric space  $\mathbb{X}$ , a *transport plan*  $\pi$  is a probability measure over  $\mathbb{X} \times \mathbb{X}$  such that for any  $A \times B \subset \mathbb{X} \times \mathbb{X}$ ,  $\pi(A \times \mathbb{X}) = \mu(A)$  and  $\pi(\mathbb{X} \times B) = \nu(B)$ . Let  $\Gamma(\mu, \nu)$  be the set of all transport plans between measures  $\mu$  and  $\nu$ . The Wasserstein distance is then defined as the minimum transport cost over  $\Gamma(\mu, \nu)$ :

$$W_2(\mu, \nu) = \sqrt{\min_{\pi \in \Gamma(\mu, \nu)} \int_{\mathbb{X} \times \mathbb{X}} d_{\mathbb{X}}(x, y)^2 d\pi(x, y)},$$

where  $d_{\mathbb{X}}(x, y)$  is the distance between  $x$  and  $y$  in the metric space  $\mathbb{X}$ . The distance to a measure is stable with respect to the Wasserstein distance as shown in [4]:

► **Theorem 4.1** (Theorem 3.5 of [4], Theorem 3.2 of [2]). *Let  $\mu$  and  $\nu$  be two probability measures on  $\mathbb{X}$  and  $m \in (0, 1]$ . Then,  $\|d_{\mu,m} - d_{\nu,m}\|_{\infty} \leq \frac{1}{\sqrt{m}} W_2(\mu, \nu)$ .*

We will mainly use the distance to empirical measures in this paper. (See [2, 4, 12] for more details on distance to a measure and its approximation.) Given a finite point set  $P$ , its associated *empirical measure*  $\mu_P$  is defined as the sum of Dirac masses:  $\mu_P = \frac{1}{|P|} \sum_{p \in P} \delta_p$ . The distance to this empirical measure for a point  $x$  can then be expressed as an average of its distances to the  $k = m|P|$  nearest neighbors where  $m$  is the mass parameter. For the sake of simplicity,  $k$  will be assumed to be an integer. The results also hold for other values of  $k$ . However, a non integer  $k$  introduces unnecessary technical difficulties. Denoting by  $p_i(x)$  the  $i$ -th nearest neighbors of  $x$  in  $P$ , one can write:

$$d_{\mu_P,m}(x) = \sqrt{\frac{1}{k} \sum_{i=1}^k d(p_i(x), x)^2}.$$

#### Geometric sampling condition

Our sampling condition treats the input point data as a measure and relates it to the manifold (where input points are sampled from) via distance-to-measures with the help of two parameters.

► **Definition 4.2.** Let  $P \subset \mathbb{R}^n$  be a discrete sample and  $M \subset \mathbb{R}^n$  a smooth manifold. Let  $\mu_P$  denote the empirical measure of  $P$ . For a fixed mass parameter  $m > 0$ , we say that  $P$  is an  $(\varepsilon, r)$ -sample of  $M$  if the following holds:

$$\forall x \in M, d_{\mu_P,m}(x) \leq \varepsilon; \quad \text{and} \tag{3}$$

$$\forall x \in \mathbb{R}^n, d_{\mu_P,m}(x) \leq r \implies d(x, M) \leq d_{\mu_P,m}(x) + \varepsilon. \tag{4}$$

The parameter  $\varepsilon$  captures the distance to the empirical measure for points in  $M$  and intuitively tells us how dense  $P$  is in relation to the manifold  $M$ . The parameter  $r$  intuitively indicates how far away we can deviate from the manifold, while keeping the noise sparse enough so as not to be mistaken for signal. We remark that if a point set is an  $(\varepsilon, r)$ -sample of  $M$  then it is an  $(\varepsilon', r')$ -sample of  $M$  for any  $\varepsilon' \geq \varepsilon$  and  $r' \leq r$ . In general, the smaller  $\varepsilon$  is and the bigger  $r$  is, the better an  $(\varepsilon, r)$ -sample is.

For convenience, denote the distance function to the manifold  $M$  by  $d_\pi : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $x \mapsto d(x, M)$ . We have the following interleaving relation:

$$\forall \alpha < r - \varepsilon, d_\pi^{-1}([-\infty, \alpha]) \subset d_{\mu_P, m}^{-1}([-\infty, \alpha + \varepsilon]) \subset d_\pi^{-1}([-\infty, \alpha + 2\varepsilon]) \tag{5}$$

To see why this interleaving relation holds, let  $x$  be a point such that  $d(x, M) \leq \alpha$ . Thus  $d(\pi(x), x) \leq \alpha$ . Using the hypothesis (3), we get that  $d_{\mu_P, m}(\pi(x)) \leq \varepsilon$ . Given that the distance to a measure is a 1-Lipschitz function we then obtain that  $d_{\mu_P, m}(x) \leq \varepsilon + \alpha$ .

Now let  $x$  be a point such that  $d_{\mu_P, m}(x) \leq \alpha + \varepsilon \leq r$ . Using the condition on  $r$  in (4) we get that  $d(x, M) \leq d_{\mu_P, m}(x) + \varepsilon \leq \alpha + 2\varepsilon$  which concludes the proof of Eqn (5).

Eqn (5) gives an interleaving between the sub-level sets of the distance to the measure  $\mu$  and the offsets of the manifold  $M$ . By Theorem 2.1, this implies the proximity between the persistence modules of their respective sub-level sets filtrations. Observe that this relation is in some sense analogous to the one obtained when two compact sets  $A$  and  $B$  have Hausdorff distance of at most  $\varepsilon$ :

$$\forall \alpha, d_A^{-1}([-\infty, \alpha]) \subset d_B^{-1}([-\infty, \alpha + \varepsilon]) \subset d_A^{-1}([-\infty, \alpha + 2\varepsilon]). \tag{6}$$

**Relation to other sampling conditions**

Our sampling condition encompasses several other existing sampling conditions. While the parameter  $\varepsilon$  is natural, the parameter  $r$  may appear to be artificial. It bounds the distances at which we can observe the manifold through the scope of the distance to a measure. In most classical sampling conditions,  $r$  is equal to  $\infty$  and thus we obtain a similar relation as for the classical Hausdorff sampling condition in Eqn (6).

One notable noise model where  $r \neq \infty$  is when there is a uniform background noise in the ambient space  $\mathbb{R}^d$ , sometimes called *clutter noise*. In this case,  $r$  depends on the difference between the density of the relevant data and the density of the noise. For other sampling conditions like Wasserstein, Gaussian, Hausdorff sampling conditions,  $r = \infty$ . Detailed relations and proofs for the Wasserstein and Gaussian sampling conditions can be found in the full version [1].

**4.2 Scalar field analysis under geometric noise**

In the rest of the paper, we assume that  $M$  is a manifold with positive reach  $\rho_M$  (minimum distance between  $M$  and its medial axis) and whose curvature is bounded by  $c_M$ . Assume that the input  $P$  is an  $(\varepsilon, r)$ -sample of  $M$  for a given  $m \in (0, 1]$ , where

$$\varepsilon \leq \frac{\rho_M}{6}, \text{ and } r > 2\varepsilon. \tag{7}$$

As discussed at the beginning of this section, we assume that there is no intrinsic functional noise, that is, for every  $p \in P$ , the observed function value  $\tilde{f}(p) = f(\pi(p))$  is the same as the true value for the projection  $\pi(p) \in M$  of this point. Our goal now is to show how to recover the persistence diagram induced by  $f : M \rightarrow \mathbb{R}$  from its observations  $\tilde{f} : P \rightarrow \mathbb{R}$  on  $P$ .

Taking advantage of the interleaving (5), we can use the distance to the empirical measure to filter the points of  $P$  to remove geometric noise. In particular, we consider the set

$$L = P \cap d_{\mu_{P,m}}^{-1}([\infty, \eta]) \text{ where } \eta \geq 2\varepsilon. \quad (8)$$

We will then use a similar approach as the one from [6] for this set  $L$ . The optimal choice for the parameter  $\eta$  is  $2\varepsilon$ . However, any value with  $\eta \leq r$  and  $\eta + \varepsilon < \rho_M$  works as long as there exist  $\delta$  and  $\delta'$  satisfying the conditions stated in Theorem 2.4.

Let  $\bar{L} = \{\pi(x) | x \in L\}$  denote the orthogonal projection of  $L$  onto  $M$ . To simulate sub-level sets  $f^{-1}([\infty, \alpha])$  of  $f : M \rightarrow \mathbb{R}$ , consider the restricted sets  $L_\alpha := L \cap (f \circ \pi)^{-1}([\infty, \alpha])$  and let  $\bar{L}_\alpha = \pi(L_\alpha)$ . By our assumption on the observed function  $\tilde{f} : P \rightarrow \mathbb{R}$ , we have:  $L_\alpha = \{x \in L | \tilde{f}(x) \leq \alpha\}$ .

Let us first recall a result about the relation between Riemannian and Euclidian metrics (e.g. [9]). For any two points  $x, y \in M$  with  $d(x, y) \leq \frac{\rho_M}{2}$  one has:

$$d(x, y) \leq d_M(x, y) \leq \left(1 + \frac{4d(x, y)^2}{3\rho_M^2}\right) d(x, y) \leq \frac{4}{3}d(x, y). \quad (9)$$

As a direct consequence of our sampling condition, for each point  $x \in M$ , there exists a point  $p \in L$  at distance less than  $2\varepsilon$ : Indeed, for each  $x \in M$ , since  $d_{\mu_{P,m}}(x) \leq \varepsilon$ , there must exist a point  $p \in P$  such that  $d(x, p) \leq \varepsilon$ . On the other hand, since the distance to measure is 1-Lipschitz, we have  $d_{\mu_{P,m}}(p) \leq d_{\mu_{P,m}}(x) + d(x, p) \leq 2\varepsilon$ . Hence  $p \in L$  as long as  $\eta \geq 2\varepsilon$ . We will use the *extrinsic* Vietoris-Rips complex built on top of points from  $L$  to infer the scalar field topology. Using the previous relation Eqn (9), we obtain the following result which states that the Euclidean distance for nearby points in  $L$  approximates the geodesic distance on  $M$ .

► **Proposition 4.3.** *Let  $\lambda = \frac{4}{3} \frac{\rho_M}{\rho_M - (\eta + \varepsilon)}$ , and assume that  $2\varepsilon \leq \eta \leq r$  and  $\varepsilon + \eta < \rho_M$ . Let  $x, y \in L$  be two points from  $L$  such that  $d(x, y) \leq \frac{\rho_M}{2} - \frac{\eta + \varepsilon}{2}$ . Then,*

$$\frac{d_M(\pi(y), \pi(x))}{\lambda} \leq d(x, y) \leq 2(\eta + \varepsilon) + d_M(\pi(x), \pi(y)).$$

**Proof.** Let  $x$  and  $y$  be two points of  $L$  such that  $d(x, y) \leq \frac{\rho_M}{2} - \frac{\eta + \varepsilon}{2}$ . As  $d_{\mu_{P,m}}(x) \leq \eta \leq r$ , Eqn (4) implies  $d(\pi(x), x) \leq \eta + \varepsilon$ . Therefore,  $d(\pi(x), \pi(y)) \leq \frac{\rho_M}{\rho_M - (\eta + \varepsilon)} d(x, y)$  [11, Theorem 4.8,(8)]. This implies  $d(\pi(x), \pi(y)) \leq \frac{\rho_M}{2}$  and following (9),  $d_M(\pi(x), \pi(y)) \leq \frac{4}{3}d(\pi(x), \pi(y))$ .

This proves the left inequality in the Proposition. The right inequality follows from

$$d(x, y) \leq d(\pi(x), x) + d(\pi(y), y) + d_M(\pi(x), \pi(y)) \leq 2(\eta + \varepsilon) + d_M(\pi(x), \pi(y)).$$

◀

► **Theorem 4.4.** *Let  $M$  be a compact Riemannian manifold and let  $f : M \rightarrow \mathbb{R}$  be a  $c$ -Lipschitz function. Let  $P$  be an  $(\varepsilon, r)$ -sample of  $M$ , and  $L$  be as introduced in Eqn (8). Assume  $\varepsilon \leq \frac{\rho_M}{6}$ ,  $r > 2\varepsilon$ , and  $2\varepsilon \leq \eta \leq r$ . Then, for any  $\delta \geq 2\eta + 6\varepsilon$  and any  $\delta' \in \left[2\eta + 2\varepsilon + \frac{8}{3} \frac{\rho_M}{\rho_M - (\eta + \varepsilon)} \delta, \frac{3}{4} \frac{\rho_M - (\eta + \varepsilon)}{\rho_M} \varrho(M)\right]$ ,  $H_*(f)$  and  $H_*(R_\delta(L_\alpha) \hookrightarrow R_{\delta'}(L_\alpha))$  are  $\frac{4}{3} \frac{c\rho_M \delta'}{\rho_M - (\eta + \varepsilon)}$ -interleaved.*

**Proof.** First, note that  $\bar{L}$  is a  $2\varepsilon$ -sample of  $M$  in its geodesic metric. It follows from the definition of  $d_{\mu_{P,m}}$  that, for any point  $x \in M$ , the nearest point  $p \in L$  to  $x$  satisfies

$d(x, p) \leq d_{\mu_P, m}(x) \leq \varepsilon$ . Hence  $d(x, \pi(p)) \leq d(x, p) + d(p, \pi(p)) \leq 2d(x, p) \leq 2\varepsilon$ . Now we apply Theorem 2.4 to  $\bar{L}$  by using  $\tilde{d}(\pi(x), \pi(y)) := d(x, y)$ ; and setting  $\lambda = \mu = \frac{4}{3} \frac{\rho_M}{\rho_M - (\eta + \varepsilon)}$ ,  $\nu = 2(\eta + \varepsilon)$ : the requirement on the distance function  $\tilde{d}$  in Theorem 2.4 is satisfied due to Proposition 4.3. The claim then follows.  $\blacktriangleleft$

Since  $M$  is compact,  $f$  is bounded due to the Lipschitz condition. We can look at the limit when  $\alpha \rightarrow \infty$ . There exists a value  $T$  such that for any  $\alpha \geq T$ ,  $L_\alpha = L$  and  $f^{-1}((-\infty, \alpha]) = M$ . The above interleaving means that  $H_*(M)$  and  $H_*(R_\delta(L)) \hookrightarrow R_{\delta'}(L)$  are interleaved. However, both objects do not depend on  $\alpha$  and this gives the following inference result:

► **Corollary 4.5.**  $H_*(M)$  and  $H_*(R_\delta(L)) \hookrightarrow R_{\delta'}(L)$  are isomorphic under conditions specified in Theorem 4.4.

## 5 Scalar Field Topology Inference under Geometric and Functional Noise

Our constructions can be combined to analyze scalar fields in a more realistic setting. Our *combined sampling condition* follows conditions (3) and (4) for the geometry. We adapt condition (2) to take into account the geometry and introduce the following conditions: we assume that there exist  $\eta \geq 2\varepsilon$  and  $s$  such that:

$$\forall p \in d_{\mu, m}^{-1}((-\infty, \eta]), |\{q \in NN_k(p) \mid |\tilde{f}(q) - f(\pi(p))| \leq s\}| \geq k' \tag{10}$$

Note that in (10), we are using  $f(\pi(p))$  as the “true” function value at a sample  $p$  which may be off the manifold  $M$ . The condition on the functional noise is only for points close to the manifold (under the distance to a measure). Combining the methods from the previous two sections, we obtain the *combined noise algorithm* where  $\eta$  is a parameter greater than  $2\varepsilon$ .

We propose the following 3-steps algorithm. It starts by handling outliers in the geometry then it makes a regression on the function values to obtain a smoothed function  $\hat{f}$  before running the existing algorithm for scalar field analysis [6] on the filtration  $\hat{L}_\alpha = \{p \in L \mid \hat{f}(p) \leq \alpha\}$ .

---

### Combined noise algorithm

1. Compute  $L = P \cap d_{\mu, m}^{-1}((-\infty, \eta])$ .
  2. Replace functional values  $\tilde{f}$  by  $\hat{f}$  for points in  $L$  using either k-median or disparity based method.
  3. Run the scalar field analysis algorithm from [6] on  $(L, \hat{f})$ .
- 

► **Theorem 5.1.** Let  $M$  be a compact smooth manifold embedded in  $\mathbb{R}^d$  and  $f$  a  $c$ -Lipschitz function on  $M$ . Let  $P \subset \mathbb{R}^d$  be a point set and  $\tilde{f} : P \rightarrow \mathbb{R}$  be observed function values such that hypotheses (3), (4), (7) and (10) are satisfied. For  $\eta \geq 2\varepsilon$ , the combined noise algorithm has the following guarantees:

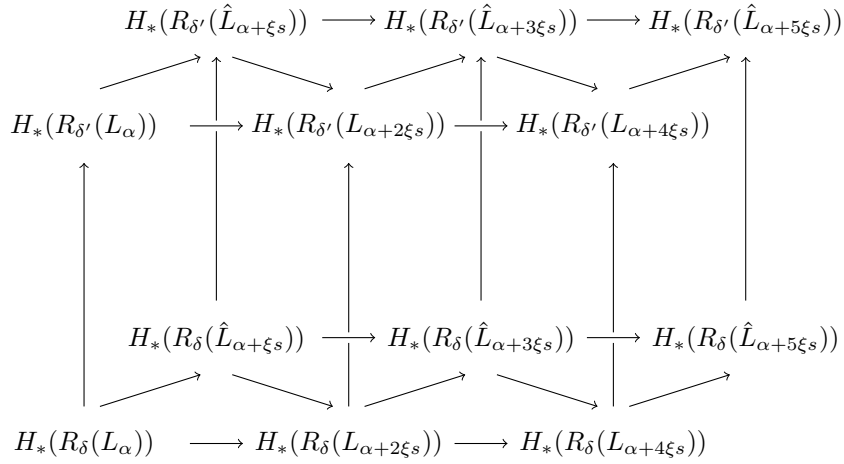
For any  $\delta \in \left[2\eta + 6\varepsilon, \frac{c(M)}{2}\right]$  and any  $\delta' \in \left[2\eta + 2\varepsilon + \frac{8}{3} \frac{\rho_M}{\rho_M - (\eta + \varepsilon)} \delta, \frac{3}{4} \frac{\rho_M - (\eta + \varepsilon)}{\rho_M} \varrho(M)\right]$ ,  $H_*(f)$  and  $H_*(R_\delta(\hat{L}_\alpha) \hookrightarrow R_{\delta'}(\hat{L}_\alpha))$  are  $\left(\frac{4}{3} \frac{c\rho_M \delta'}{\rho_M - (\eta + \varepsilon)} + \xi s\right)$ -interleaved where  $\xi = 1$  if we use the  $k$ -median and  $\xi = \left(1 + 2\sqrt{\frac{k - k'}{2k' - k}}\right)$  if we use the disparity method for Step 2.

**Proof.** First, consider the filtration induced by  $L_\alpha = \{x \in L | f(\pi(x)) \leq \alpha\}$ ; that is, we first imagine that all points in  $L$  have correct function values (equals to the true value of their projection on  $M$ ). By Theorem 4.4, for

$$\delta \in \left[ 2\eta + 6\varepsilon, \frac{\varrho(M)}{2} \right] \text{ and } \delta' \in \left[ 2\eta + 2\varepsilon + \frac{8}{3} \frac{\rho_M}{\rho_M - (\eta + \varepsilon)} \delta, \frac{3}{4} \frac{\rho_M - (\eta + \varepsilon)}{\rho_M} \varrho(M) \right],$$

$H_*(f)$  and  $H_*(R_\delta(L_\alpha) \hookrightarrow R_{\delta'}(L_\alpha))$  are  $\frac{4}{3} \frac{c\rho_M\delta'}{\rho_M - (\eta + \varepsilon)}$ -interleaved.

Next, consider  $\hat{L}_\alpha = \{p \in L | \hat{f}(p) \leq \alpha\}$ , which leads to a filtration based on the smoothed function values  $\hat{f}$  (not observed values). Recall that our algorithm returns  $H_*(R_\delta(\hat{L}_\alpha) \hookrightarrow R_{\delta'}(\hat{L}_\alpha))$ . We aim to relate this persistence module with  $H_*(R_\delta(L_\alpha) \hookrightarrow R_{\delta'}(L_\alpha))$ . Specifically, fix  $\alpha$  and let  $(x, y)$  be an edge of  $R_\delta(L_\alpha)$ . This means that  $d(x, y) \leq 2\delta$ ,  $f(\pi(x)) \leq \alpha$ ,  $f(\pi(y)) \leq \alpha$ . Corollary 3.2 can be applied to the function  $f \circ \pi$  due to hypothesis (10). Hence  $|\hat{f}(x) - f(\pi(x))| \leq \xi s$  and  $|\hat{f}(y) - f(\pi(y))| \leq \xi s$ . Thus  $(x, y) \in R_\delta(\hat{L}_{\alpha + \xi s})$ . One can reverse the role of  $\hat{f}$  and  $f$  and get an  $\xi s$ -interleaving of  $\{R_\delta(L_\alpha)\}$  and  $\{R_\delta(\hat{L}_\alpha)\}$ . This gives rise to the following commutative diagram since all arrows are induced by inclusions.



Thus the two persistence modules induced by filtrations of nested pairs  $\{R_\delta(L_\alpha) \hookrightarrow R_{\delta'}(L_\alpha)\}$  and  $\{R_\delta(\hat{L}_\alpha) \hookrightarrow R_{\delta'}(\hat{L}_\alpha)\}$  are  $\xi s$ -interleaved. Combining this with the interleaving between  $H_*(R_\delta(L_\alpha) \hookrightarrow R_{\delta'}(L_\alpha))$  and  $H_*(f)$ , we obtain the stated results. ◀

We note that, while this theorem assumes a setting where we can ensure theoretical guarantees, the algorithm can be applied in a more general setting still producing good results.

**Acknowledgments.** This work was supported by the ANR project TopData 13-BS01-008, the ERC project Gudhi 339025 and the NSF grants CCF-1064416, CCF-1116258, CCF-1319406 and CCF-1318595.

References

- 1 M. Buchet, F. Chazal, T. K. Dey, F. Fan, S. Y. Oudot, and Y. Wang. Topological analysis of scalar fields with outliers. *arXiv preprint arXiv:1412.1680*, 2014.
- 2 M. Buchet, F. Chazal, S. Oudot, and D. R. Sheehy. Efficient and robust persistent homology for measures. In *Proceedings of the 26th ACM-SIAM symposium on Discrete algorithms*. SIAM, 2015.

- 3 F. Chazal, D. Cohen-Steiner, M. Glisse, L. J. Guibas, and S. Oudot. Proximity of persistence modules and their diagrams. In *Proc. 25th ACM Sympos. on Comput. Geom.*, pages 237–246, 2009.
- 4 F. Chazal, D. Cohen-Steiner, and Q. Mérigot. Geometric inference for probability measures. *Foundations of Computational Mathematics*, 11(6):733–751, 2011.
- 5 F. Chazal, V. de Silva, M. Glisse, and S. Oudot. The structure and stability of persistence modules, 2013. arXiv:1207.3674.
- 6 F. Chazal, L. J. Guibas, S. Y. Oudot, and P. Skraba. Scalar field analysis over point cloud data. *Discrete & Computational Geometry*, 46(4):743–775, 2011.
- 7 F. Chazal and S. Y. Oudot. Towards persistence-based reconstruction in euclidean spaces. In *Proceedings of the twenty-fourth annual symposium on Computational geometry*, pages 232–241. ACM, 2008.
- 8 D. Cohen-Steiner, H. Edelsbrunner, and J. Harer. Stability of persistence diagrams. *Discrete & Computational Geometry*, 37(1):103–120, 2007.
- 9 T. K. Dey, J. Sun, and Y. Wang. Approximating cycles in a shortest basis of the first homology group from point data. *Inverse Problems*, 27(12):124004, 2011.
- 10 H. Edelsbrunner and J. Harer. *Computational Topology: An Introduction*. Amer. Math. Soc., Providence, Rhode Island, 2009.
- 11 H. Federer. Curvature measures. *Transactions of the American Mathematical Society*, pages 418–491, 1959.
- 12 L. Guibas, D. Morozov, and Q. Mérigot. Witnessed k-distance. *Discrete & Computational Geometry*, 49(1):22–45, 2013.
- 13 L. Györfi. *A distribution-free theory of nonparametric regression*. Springer, 2002.
- 14 J. Kloke and G. Carlsson. Topological de-noising: Strengthening the topological signal. *arXiv preprint arXiv:0910.5947*, 2009.
- 15 S. Kpotufe. k-nn regression adapts to local intrinsic dimension. *arXiv preprint arXiv:1110.4300*, 2011.
- 16 A. Zomorodian and G. Carlsson. Computing persistent homology. *Discrete & Computational Geometry*, 33(2):249–274, 2005.