Stabilizing the unstable output of persistent homology computations

Paul Bendich, Peter Bubenik, and Alexander Wagner

ABSTRACT. We propose a general technique for extracting a larger set of stable information from persistent homology computations than is currently done. The persistent homology algorithm is usually viewed as a procedure which starts with a filtered complex and ends with a persistence diagram. This procedure is stable (at least to certain types of perturbations of the input). This justifies the use of the diagram as a signature of the input, and the use of features derived from it in statistics and machine learning. However, these computations also produce other information of great interest to practitioners that is unfortunately unstable. For example, each point in the diagram corresponds to a simplex whose addition in the filtration results in the birth of the corresponding persistent homology class, but this correspondence is unstable. In addition, the persistence diagram is not stable with respect to other procedures that are employed in practice, such as thresholding a point cloud by density. We recast these problems as real-valued functions which are discontinuous but measurable, and then observe that convolving such a function with a suitable function produces a Lipschitz function. The resulting stable function can be estimated by perturbing the input and averaging the output. We illustrate this approach with a number of examples, including a stable localization of a persistent homology generator from brain imaging data.

1. Introduction

Persistence diagrams, also called bar codes, are one of the main tools in topological data analysis [10, 35, 31, 37]. In combination with machine-learning and statistical techniques, they have been used in a wide variety of real-world applications, including the assessment of road network reconstruction [3], neuroscience [23], [6], vehicle tracking [5], object recognition [39], protein compressibility [36], and protein structure [38].

Put briefly, these persistence diagrams are multi-sets of points in the extended plane, and they compactly describe some of the multi-scale topological and geometric information present in a high-dimensional point cloud, or carried by a real-valued function on a domain. Several theorems [24, 14, 26] state that persistence diagrams are stable with respect to certain variations in the point-cloud or functional input, and so the conclusions drawn from them can be taken with some confidence.

1

On the other hand, there is additional potentially very useful but unstable information produced during the computation of persistence diagrams. For example, a point far from the diagonal in the degree-zero persistence diagram represents a connected component with high persistence. This component first appears somewhere and the computation that produces the persistence diagram can be used to find its location. However this location is not stable: as we will describe below, a small change in the input will cause only a small change in the persistence of this connected component, but it can radically alter the location of its birth. We summarize this as follows.

Fundamental Conundrum of Topological Data Analysis. Users of topological data analysis would like to find the simplices or cycles corresponding to the birth of the most significant pairings of critical values. However, unlike the paired critical values, these simplices and cycles are unstable.

In addition, persistent homology computations may rely on parameters such that the output persistence diagram is not stable with respect to changes of these parameters.

1.1. Our Contribution. This paper introduces a method for stabilizing desirable but unstable outputs of persistent homology computations. The main idea is the following. On the front end, we think of a persistent homology computation C as being parametrized by a vector $\alpha = (\alpha_1, \dots, \alpha_n)$ of real numbers. These parameters could specify the input to the computation (e.g. the coordinates of the vertices of a simplicial complex) or they could specify other values used in the computation (e.g. threshold parameters used in de-noising or bandwidths for smoothing). For a given choice of a, we get a persistence diagram. On the back end, we consider a function p that extracts a real-number summary from a persistence diagram. For example, p might extract the persistence of a homology class created by the addition of a specific edge in a filtered simplicial complex, or it might be an indicator function on whether or not the longest bar was born by the addition of a simplex contained in a fixed region of the input space. The composite function h that maps the parameter vector to the real number need not be continuous, but it will in many cases be measurable. We convolve this function with a Gaussian (or indeed any Lipschitz function) to produce a new Lipschitz function that carries the persistence-based information we desire.

Our main theoretical results (Theorems 4.2, 4.3, and 4.4) give conditions on functions h and K (where K will usually be a kernel) that guarantee that the convolution h * K is Lipschitz with specified Lipschitz constant. From these we obtain the following, where more precise statements are given as Corollaries 4.5, 4.6, and 4.8.

Theorem 1.1. If h is locally essentially bounded then for the triangular and Epanechnikov kernels, h * K is locally Lipschitz. If h is essentially bounded then for the Gaussian kernel, h * K is Lipschitz.

In practice, this can be translated to a simple procedure for stabilizing unstable persistent homology computations: perturb the input by adding, for example, Gaussian noise, and redo the computation; repeat and average. By the law of large numbers, the result converges to the desired stable value.

THEOREM 1.2. Let $\epsilon_1, \ldots, \epsilon_M$ be drawn independently from a kernel K. Then

$$\frac{1}{M}\sum_{i=1}^{M}h(\alpha-\varepsilon_{i})\rightarrow (h*K)(\alpha).$$

We summarize our computational pipeline in the following algorithm. Say we have performed a persistence computation and obtained an unstable output. For example, we have determined that the longest interval in the degree one bar code of the Vietoris-Rips complex on points $X_1, \ldots, X_N \in \mathbb{R}^d$ is born with the addition of the edge X_1X_2 . We encode this output as a function $h: \mathbb{R}^n \to \mathbb{R}$ with input $a \in \mathbb{R}^n$. For example, the coordinates of the above points give us $a \in \mathbb{R}^n$ where n = Nd. We define a function $h: \mathbb{R}^n \to \mathbb{R}$ whose value is the length of the longest interval in the bar code if it is born with the addition of the edge X_1X_2 and is otherwise 0.

Algorithm 1 Stabilizing unstable persistence computations

```
Input: h: \mathbb{R}^n \to \mathbb{R}, a \in \mathbb{R}^n

Parameters: M \in \mathbb{N} <, \sigma > 0

for i \leftarrow 1, M do

for j \leftarrow 1, n do

Sample \varepsilon_j from N(0, \sigma^2)

end for

y_i \leftarrow h(a + \varepsilon), \varepsilon = (\varepsilon_1, \dots, \varepsilon_n)

end for

return the average value of y_1, \dots, y_M
```

The choice of standard deviation σ is discussed in Section 4.6. In Section 4.5, we prove that Algorithm 1 is stable with respect to this choice.

- **1.2. Two examples.** For the reader familiar with persistent homology who wants to see how this works in practice, we provide two examples, the first to real data and the second to synthetic data. Code for these examples is available at https://github.com/peter-bubenik/stabilizing-paper-code.
- 1.2.1. Location of a persistent homology generator. In [6], the authors apply topological data analysis to brain arteries extracted from magnetic resonance images. Mathematically, each of these brain arteries is a graph embedded in three-dimensional Euclidean space. Using the height (the z-coordinate) one obtains a filtration on this graph, which may be used to compute degree-zero persistent homology.

To facilitate statistical analysis of the resulting persistence diagrams (i.e. bar codes), they convert each persistence diagram to a vector consisting of the lengths of the 100 longest bars in decreasing order. In their analysis, the length of the 28th longest bar is a numerical feature that yields a correlation with age that is near-optimal among vector features consisting of the lengths of the ith through jth longest bars for any $1 \le i \le j \le 100$.

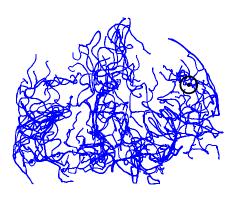


FIGURE 1. The brain arteries of the first subject in [6]. The black dot is the location of the generator of the 28th longest bar in degree-zero persistent homology. We consider the indicator function on the location of this generator with respect to the given sphere.

If one wants to find a biological interpretation of this result, it is obvious to ask for the location of the generator of the 28th longest bar for each subject. It is easy to locate the generator responsible for the birth of the 28th longest bar. It will be a particular vertex of the graph, whose image is a point in space. However, the location of this point is unstable: as we will later explain, small perturbations of the spatial coordinates of the vertices of the graph can lead to large changes of this location.

We choose a ball centered at this point and consider the function whose value is 1 if the location of the generator of the 28th longest bar is located in this ball, and is otherwise 0. The resulting function $h: \mathbb{R}^{3V} \to \mathbb{R}$ (where V is the number of vertices in the graph) is unstable, but it may be stabilized using the method summarized in Section 1.1. Applying Algorithm 1 with M=1000 and $\sigma=0.1$ we obtain an estimate of the stable value of h*K evaluated at the observed input, equal to 0.471. This shows that under small perturbations of the input, roughly half of the time is the generator of the 28th longest bar located in the chosen ball. This result holds for a large range of sizes of balls - see Appendix A for some further discussion.

We remark that this approach provides a resolution of the conflict between TDA theorists and TDA users expressed in the Fundamental Conundrum of Topological Data Analysis 1. We can provide TDA users with a location of a generator of a persistent homology class together with an estimate of a stable real value of how often this location lies in a given region under certain perturbations.

1.2.2. Persistence of a homology class born in a region. Consider the function f on the square in Figure 2. This induces a function \bar{f} on the torus since f(x,y)=0 on the boundary of the square. Suppose we are only given a finite sample of this induced function and we are interested in the presence of long-lived bars which are born in the region of the torus corresponding to the second quadrant of the square.

To be concrete, we start with a sample X of N points from the graph of \bar{f} , by sampling u_i, v_i independently from the uniform distribution on $[-\pi, \pi]$ and letting $z_i = f(u_i, v_i)$. Note that X is a random variable. We use X to construct a filtered simplicial complex

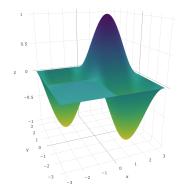
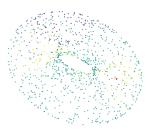


Figure 2. The graph of the function on the square $[-\pi,\pi]^2$ given by $f(\mathfrak{u},\nu)=\sin(\mathfrak{u})\sin(\nu)(1-0.9^*1_{\{\mathfrak{u}<0,\nu<0\}}):[-\pi,\pi]^2\to\mathbb{R}.$ It induces a function on the torus, $\bar{f}:\mathsf{T}^2\to\mathbb{R}$, with two global minima with value -1, one global maximum with value 1, one local maximum with value 0.1, and four saddle points with value 0. From





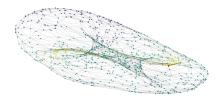


FIGURE 3. Sample of 1000 points from the graph $\{(x, f(x)) : x \in T^2\}$, where the function values are indicated using the same color scale as in Figure 2. The points on the torus are used to construct a Delaunay triangulation, which is filtered using the function values. On the right we indicate the filtration values by moving the points in the normal direction.

approximating the unknown function \bar{f} as follows. From the points $\{(u_i, v_i)\}$ we construct a Delaunay triangulation of the torus. We filter this simplicial complex by assigning the vertex (u_i, v_i) the value z_i and assigning edges and triangles the maximum value of their vertices.

We compute the 0-dimensional extended persistence diagram of this filtered simplicial complex. Let h(X) be the length of the longest bar if that bar was born in the region corresponding to the second quadrant (see Figure 2) and 0 otherwise.

This process defines a function $h : \mathbb{R}^{3N} \to \mathbb{R}$, but h is unstable. Consider the sample X = x in Figure 3. We have h(x) = 0 since the global minimum, highlighted in red, is born outside the region corresponding to the second quadrant. Because of the

¹Extended persistent homology follows the homology of increasing sublevel sets with the relative homology of the whole space relative to decreasing superlevel sets [25]. In the case considered here, it pairs the global minimum with the global maximum.

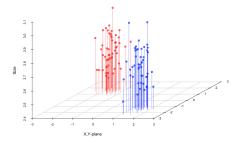


Figure 4. Locations and sizes of 100 longest bars from the trials. Averaging the lengths of the blue bars over 1000 trials we get 1.3644, which is consistent with the fact that the random variable h(X) is 0 or about 2 with equal probability. We should not expect $\lim_{M\to\infty}\frac{1}{M}\sum_{i=1}^M h(x+\varepsilon_i)$ to converge to 1 because unlike f, a particular sample X=x is not symmetric with respect to the second and fourth quadrants.

symmetry of f, the random variable h(X) is 0 approximately half the time and about 2 approximately half the time.

Let K denote the 3N-variate Gaussian with mean 0 and standard deviation 0.2. For $M \ge 1$, sample $\epsilon_1, \ldots, \epsilon_M$ independently from K. Compute $\frac{1}{M} \sum_{i=1}^M h(x - \epsilon_i)$. See Figure 4. As M increases, this quantity converges to g(x), where g := h * K is the stabilized version of h.

1.3. Related Work. Partial inspiration for the main idea of our work (when faced with an instability caused by a near-interchange of values, perturb the values many times and take some sort of average) comes from the trembling-hand equilibrium solution [40] to the non-uniqueness problem for Fréchet means of persistence diagrams. Our approach should also be compared with the topological reconstruction results of Niyogi, Smale and Weinberger [43].

Several recent papers have advocated principled approaches for extracting features from persistence diagrams, including persistence landscapes [8], the stable multi-scale kernel [45], intensity functionals [21], persistence images [2], the stable topological signature [13], and the cover-tree entropy reduction [47]. Our result complements these ideas: once one identifies some specific parts of the persistence diagram as having good classification power, one can then attempt to locate, in a robust way, the portions of the domain responsible for these parts. Other papers (e.g. [15, 9, 1]) have developed sophisticated schemes for data-cleaning before persistent homology computation. These techniques are generally fragile to certain initial parameter choices, such as the m_0 parameter in [15]. Again, we provide a complementary role: any of these schemes can be run many times for several perturbations of an initial parameter choice, and the output can then be taken with confidence.

Dey and Wenger [29] have shown that the critical points of interval persistent homology are stable in the sense that they remain within some path-connected component.

Zomorodian and Carlsson [50] use Mayer-Vietoris as inspiration in their technique for localizing (relative to a cover) homology classes within a given simplicial complex. However, this works only for a fixed simplicial complex, not a simplicial complex endowed with a filtration, and the results are certainly fragile to changes in this fixed complex.

Weinberger [49] considers the sample complexity of some basic problems of topological inference. Specifically, he estimates the number of sample points necessary to determine the dimension, topological type, and to detect singularities for certain spaces.

Robust summaries of persistent homology are considered in the following papers; they do not consider the location of homology generators. In [7], Blumberg et al. show that persistent homology on a metric measure space induces a stable empirical measure in the space of persistence diagrams. Taking the distance to a reference distribution or a reference barcode, they obtain robust statistics. In [17], the authors derive limiting distributions and confidence sets for persistence diagrams based on the sub-level sets of the distance-to-a-measure.

Convolving with a kernel to obtain smoothness is a classical idea in statistics [46, 48]. It has been used to construct smooth estimators of discrete data as an initial step to computing persistent homology [9, 8, 33]. A related idea is the to perform subsampling (e.g. the bootstrap) to obtain convergence results and confidence intervals for persistence diagrams and persistence landscapes [33, 20, 18, 19]. These papers use ideas related to ones presented here, but to smooth initial data or to smooth stable outputs of persistence computations, not to stabilize unstable outputs of persistence computations.

1.4. Outline. Persistent homology computations and stability theorems are reviewed in Section 2, although we assume the reader is already somewhat familiar with them. Several examples of important but unstable persistence-based information are given in Section 3, and we then describe a general approach that stabilizes them in Section 4. In Section 5, we show how to apply these results to persistent homology computations. Computational experiments, and some suggested interpretations of the results, are presented in Section 6. In Section 1.2.1, we apply to our methods to an application with real data. Potential future directions are discussed in Section 7.

2. Persistent Homology and Stability

The treatment of persistence diagrams here is adapted from [31]. For a more general discussion, see [44]. We assume the reader is familiar with the basics of homology groups: the textbook [41] is a good introduction. All homology groups are assumed to be computed over some fixed field. For concreteness, we restrict our attention to simplicial complexes, but our results also apply to more general complexes.

2.1. Persistent Homology.

- 2.1.1. Filtered simplicial complexes. Persistent homology is computed for a finite filtered abstract simplicial complex. That is, we have a finite abstract simplicial complex, a collection, $K = \{\sigma\}$, of nonempty subsets of a fixed finite set that satisfy the condition that if $\emptyset \neq \tau \subseteq \sigma \in K$ then $\tau \in K$. In addition, we have a filtration, a function $f: K \to \mathbb{R}$ such that $\tau \subseteq \sigma$ then $f(\tau) \leqslant f(\sigma)$. That is, f is order preserving.
- 2.1.2. Persistence Diagrams. Fix a homological dimension p. Suppose the distinct values of f are $r_1 < \ldots < r_m$. For each $1 \leqslant i \leqslant m$, define $K^i = \{\sigma \in K \mid f(\sigma) \leqslant r_i\}$. Since f is order preserving, each K^i is a subcomplex. Whenever $i \leqslant j$, there is an inclusion $K^i \hookrightarrow K^j$, which induces a homomorphism:

$$f_p^{i,j}: H_p(K^i) \to H_p(K^j).$$

A homology class $\alpha \in H_p(K^i)$ is a *persistent homology class* that is *born* at level i if $\alpha \notin \inf_p^{i-1,i}$, and that *dies* entering level j if $f_p^{i,j}(\alpha) = 0$ but $f_p^{i,j-1}(\alpha) \neq 0$. If α never dies, we say that it dies entering level $j = \infty$ and $r_\infty = \infty$. The *persistence* of α is defined to be $pers(\alpha) = r_j - r_i$. The set of classes which are born at i and die entering level j form a vector space, with rank denoted $\mu_p^{i,j}$. The degree-p *persistence diagram* of f, $pgm_p(f)$, encodes these ranks. It is a multiset of points in the extended plane, with a point of multiplicity $\mu_p^{i,j}$ at each point (r_i, r_j) .

2.1.3. *Persistent homology computations*. In practice, one constructs a filtered abstract simplicial complex from some other starting data. In addition, more information can be extracted from the persistent homology algorithm than just the persistence diagram.

We define a *persistent homology computation*, \mathfrak{C} , to be a function whose input consists of real numbers a_1, a_2, \ldots, a_n . These may include input values and also parameter values for the computation. Using this input, \mathfrak{C} constructs an abstract simplicial complex K together with a filtration f. The output of \mathfrak{C} consists of a degree-p persistence diagram together with for each (r_i, r_j) in the persistence diagram (counted with multiplicity), a p-simplex σ with $f(\sigma) = r_i$, a p-cycle α in K^i containing σ , a (p+1)-simplex τ with $f(\tau) = r_j$, and a (p+1)-chain β in K^j containing τ with $d\beta = \alpha$.

2.1.4. Examples.

Example 2.1. Functions on simplicial complexes. A filtered abstract simplicial complex, K, may be obtained from a real-valued function, F, on the vertices in a finite simplicial complex, \mathcal{K} . As a set $K \cong \mathcal{K}$. A filtration, f, on K is defined by $f(\sigma) = \sup_{x \in \sigma} F(x)$.

For example, let \mathcal{K} be the geometric line graph (i.e. an embedding of a graph consisting of vertices and edges - in the plane), shown on the bottom of the left side of Figure 5. Above this, we have the graph of a function F on the points in \mathcal{K} . From this, we have a corresponding abstract simplicial complex K and filtration f. The persistence diagram $Dgm_0(f)$ is on the right. The input to \mathcal{C} consists of the function values (from left to right) a_1, a_2, \ldots, a_n .

Example 2.2. Distance to a PL-Curve. Consider the piecewise-linear curve C on the left side of Figure 6. Moving clockwise, we order its vertices $A = v_1, v_2, \dots v_N = D$. Let K be the full simplex on these N vertices. For each vertex v, define f(v) = 0. For each

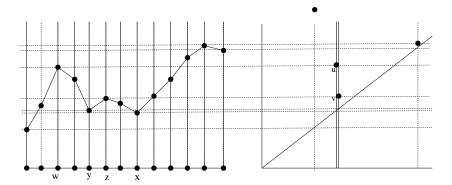


FIGURE 5. Left: The graph of a function F on a simplicial complex \mathcal{K} . Right: the degree-zero persistence diagram $\mathrm{Dgm}_0(f)$ for the corresponding abstract simplicial complex K and filtration f. The labeled points have coordinates $\mathfrak{u}=(f(x),f(w))$ and $\nu=(f(y),f(z))$. The point on the very top has infinite y-coordinate.

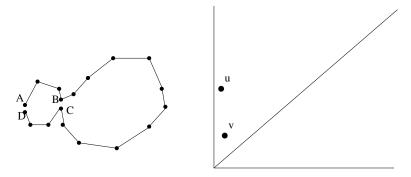


FIGURE 6. Left: a piecewise-linear curve in the plane. The distance between A and D is slightly smaller than the distance between B and C. Right: $Dgm_1(f)$, where f is as defined in the text. The points $\mathfrak u$ and $\mathfrak v$ correspond to one-cycles that are created by the additions of edges (A,D) and (B,C), respectively.

edge of the form $e = (\nu_i, \nu_{i+1})$, define f(e) = 0, and for any other edge $e = (\nu_i, \nu_j)$, we set f(e) to be the Euclidean distance between ν_i and ν_j . Finally, for any higher simplex σ , set $f(\sigma) = \max_{e \subseteq \sigma} f(e)$, where we take the maximum over the set of edges contained in σ . The degree-one persistence diagram $Dgm_1(f)$ appears on the right of Figure 6.

Here the input to C consists of the 2n coordinates of the vertices. We note this paradigm can be extended to curves C in higher-dimensional ambient spaces, or even to higher-dimensional complexes.

Example 2.3. Point cloud – Vietoris-Rips. Suppose that $X = \{x_1, \dots, x_N\}$ is a set of points in some metric space (Y, d). We let K be the full simplex on these vertices. Define f(v) = 0 for each vertex and f(e) = d(v, w) for each edge e = (v, w). As above, we set $f(\sigma) = \max_{e \subseteq \sigma} f(e)$ for all higher-dimensional simplices. This is called the Vietoris-Rips filtration. We denote $Dgm_p(X) = Dgm_p(f)$. The input to $\mathfrak C$ consists of the coordinates of

the points in X in some parametrization of Y. Alternatively, it consists of the entries of the distance matrix $D = (d(x_i, x_j))$.

For example, let X be the circular point cloud on the top-left of Figure 8. The corresponding $Dgm_1(X)$ appears on the top-right of the same figure.

Example 2.4. *Point cloud – geometry.* Often, the simplicial complex in the previous example is too large to work with. Instead one applies some geometric ideas to construct a smaller filtered simplicial complex. Examples include witness complexes [28], the graph-induced complex [30], and the use of nudged elastic bands [1]. These constructions typically include one or more parameters, which we append to the input to $\mathfrak C$.

EXAMPLE 2.5. *Point cloud – statistics*. Instead of using geometric ideas to construct a smaller point cloud we can use statistical ideas. For example, one can use a kernel to smooth the point cloud to obtain a density estimator on the underlying space Y and use this to filter a triangulation of Y [23, 9, 22]. Or one may use the local density to threshold the point cloud [11]; we consider this in more detail in Example 5.4. Again, these constructions include one or more parameters, which we append to the input for \mathcal{C} .

Example 2.6. Regression. Here we present a variant of Example 2.1 in which we are not given the simplicial complex. Instead we sample points $X=(x_1,\ldots,x_N),\,x_i\in\mathbb{R}^d$ from some probability distribution on \mathbb{R}^d . We also sample corresponding perturbed function values $y_i\in\mathbb{R}$. For example, we may have $y_i=f(x_i)+\varepsilon_i$, where ε_i is sampled from a univariate Gaussian. We use X to construct a Delaunay triangulation K. We then use $Y=(y_1,\ldots,y_N)$ to filter K as follows: $f(\sigma)=\max_{x_i\in\sigma}y_i$. This is called the lower star filtration.

Instead of the sample points lying in \mathbb{R}^d , they may lie on some compact Riemannian manifold. See the torus example in Section 1.

Instead of filtering K directly using Y, one can instead use (X,Y) to construct an estimator \hat{f} of the unknown regression function f. We can then use \hat{f} to filter K [9].

2.2. Stability. The persistence diagram $Dgm_p(f)$ is a summary of the function f, and it turns out to be a stable one. The discussion here is adapted from [24]. For a broader description, see [14, 16].

For convenience, to each persistence diagram, we add every point (r,r) on the major diagonal, each with infinite multiplicity.

Now suppose that $\phi:D\to D'$ is some bijection between two persistence diagrams; bijections exist because of the infinite-multiplicity points along the diagonal. The cost of ϕ is defined to be $C(\phi)=\sup_{u\in D}\|u-\phi(u)\|_{\infty}$; that is, the largest box-norm distance between matched points. The bottleneck distance $W_{\infty}(D,D')$ is defined to be the minimum cost amongst all such bijections. For example, if D and D' are the black and red diagrams, respectively, on the right side of Figure 7, then the best bijection would pair u with u', v with v', the two infinite-persistence points with each other,

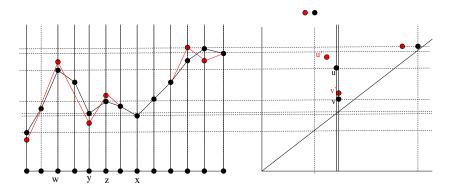


FIGURE 7. Left: The graphs of functions f (black) and g (red), both on the same domain \mathcal{K} . Right: the persistence diagrams $\mathrm{Dgm}_0(f)$ and $\mathrm{Dgm}_0(g)$, using the same color scheme.

and the other two points with the closest diagonal points. The bottleneck distance is the cost of this bijection. The Diagram Stability Theorem [24] guarantees that persistence diagrams of nearby functions are close to one another. More precisely, we have $W_{\infty}(D_{\mathfrak{p}}(f),D_{\mathfrak{p}}(g))\leqslant ||f-g||_{\infty}$. This is illustrated by Figure 7.

Note that the difference between f and g is measured in the L_{∞} norm. In the point cloud context (Ex. 2.3), this translates into requiring that the two point cloud inputs be Hausdorff-close. However, the persistence diagram is not stable with respect to the addition of outliers. We discuss this problem in more detail in Section 3.2 and propose a solution in Section 4.

3. Instability

The Diagram Stability Theorem tells us that the persistence diagram obtained in the output of a persistent homology computation is stable with respect to certain perturbations of the input used to construct a filtered abstract simplicial complex. However, other outputs of persistent homology computations are not stable. This includes the simplices and cycles that generate persistent homology classes. These are of great interest to practitioners hoping to interpret persistence calculations more directly. In addition, many persistence computations rely on choices of parameters and the resulting persistence diagrams may be unstable with respect to these choices.

3.1. Instability of Generating Cycles/Simplices. Persistence diagrams are useful and robust measures of the *size* of topological features. What they are less good at, on the other hand, is robustly pinpointing the *location* of important topological features. We use Figure **7** to illustrate this problem. Suppose that we have the fixed domain K and we observe the function f. One of the most prominent points in $Dgm_0(f)$ is u, which corresponds to the pair of values f(x) and f(w). We might thus be tempted to say that f has an important feature, a component of high-persistence, at x. But consider the nearby function g instead. Its diagram $Dgm_0(g)$ has a point u' that is very close to u, but this point corresponds to the pair of values f(y) and f(w). There is still a component born

at g(x), but it corresponds to the much smaller persistence point v'. And so while the persistence of the point u is a stable summary of the function f, the actual location x of the topological feature it corresponds to is not.

This is unfortunate. Several recent works ([5], [6], among others) have shown that the presence of points in certain regions of the persistence diagram has strong correlation with covariates under study. For example, each diagram in the second cited work came from a filtration of the brain artery tree in a specific patient's brain, and it was found that the density of points in a certain middle-persistence range gave strong correlations with patient age. It would of course be tempting to hold specific locations in the brain responsible for these points with high distinguishing power.

Unsurprisingly, this problem remains for persistent homology in higher degrees. Consider Figure 6 again. It is easy to see that edge (A,D) creates the large loop which corresponds to point $u \in Dgm_1(f)$. However, a slight perturbation of the vertex configuration could render (B,C) responsible for this loop instead, and so we cannot robustly locate the persistence of this loop at(A,D).

In Section 4, we both rigorously define this non-robustness and suggest a method for addressing it.

3.2. Outliers and Instability of Parameter Choices. The Diagram Stability Theorem guarantees the persistence diagrams associated to two Hausdorff-close point clouds will themselves be close. However, it says nothing about the outlier problem. For example, consider again the point cloud X (Figure 8, top-left) from Example 2.3 to which we apply the Vietoris-Rips construction. Its persistence diagram $Dgm_1(X)$ (top-right of same figure) has one high-persistence point, which corresponds to the "circle" that we qualitatively see when looking at the points. On the other hand, consider the point cloud X' on the bottom-left, which consists of X and three "outlier" points spread across the interior of the circle. The diagram $Dgm_1(X')$ (bottom-right) is not close to $Dgm_1(X)$: there is still one point of fairly high persistence, but it's much closer to the diagonal than before.

In practice, this problem is often addressed by first de-noising the point cloud in some way. For example, Carlsson et. al. [12] first thresholded by density before computing Vietoris-Rips filtrations when they discovered a Klein bottle in the space of natural images. There are no guarantees that a different, nearby choice of density threshold parameter would not give a qualitatively different persistence diagram. Section 4 addresses this by introducing a general method for handling parameter choice in persistence computations.

4. Theory: Stability from convolutions

In this section we show how functions may be stabilized by convolving them with a kernel. In Section 5, we will apply these results to the function $h:\mathbb{R}^n\to\mathbb{R}$ discussed in the introduction. First, we give three general results with various assumptions on the function and the kernel. Next, we apply them in three particular cases: the simple

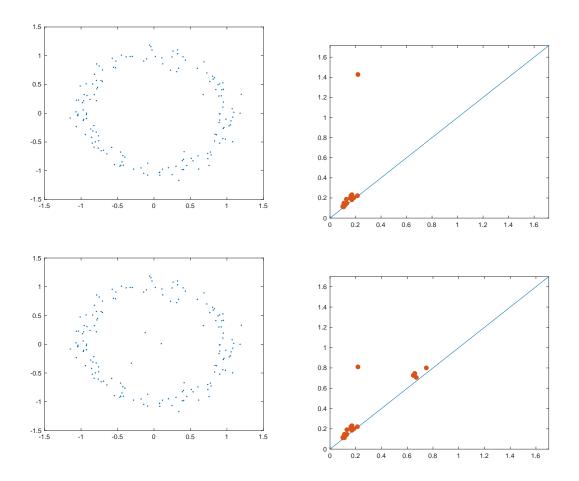


FIGURE 8. Illustration of the outlier problem for the persistent homology of the Vietoris-Rips complex of a point cloud. All figures produced in MATLAB. Top left: 150 points X, sampled with bounded noise from a circle. Top right: $Dgm_1(X)$. Bottom left: 153 points Y, which is X plus three outlier points. Bottom right: $Dgm_1(Y)$.

triangular kernel and the commonly used Epanechnikov and Gaussian kernels. We then outline a few specific examples, some of which will be explored via experiment in the next section.

4.1. Lipschitz functions and convolution. Let us start by recalling a few definitions. For $C \ge 0$, a function $f: \mathbb{R}^n \to \mathbb{R}$ is said to be *C-Lipschitz* if for all $\mathfrak{u}, \mathfrak{v} \in \mathbb{R}^n$, $|f(\mathfrak{u}) - f(\mathfrak{v})| \le C|\mathfrak{u} - \mathfrak{v}|$, where $|\mathfrak{x}|$ denotes the Euclidean norm. We will call a function *Lipschitz* if it is C-Lipschitz for some $C \ge 0$. The support of f, denoted supp(f), is the closure of the subset of \mathbb{R}^n where f is non-zero.

Let $h, g : \mathbb{R}^n \to \mathbb{R}$ be (Lebesgue) measurable functions that are defined almost everywhere. The *1-norm* of h, is given by $\|h\|_1 = \int_{\mathbb{R}^n} |h(t)| dt$, if it exists. The *essential supremum* of h, denoted by $\|h\|_{\infty}$, is the smallest number a such that the set $\{x \mid |f(x)| > a\}$ has

measure 0. If it exists, the *convolution product* of h and g, is given by

$$(h*g)(t) = \int_{\mathbb{R}^n} h(s)g(t-s)ds = \int_{\mathbb{R}^n} h(t-s)g(s)ds.$$

It exists everywhere, for example, if one function is essentially bounded and the other is integrable; or if one function is bounded and compactly supported and the other is locally integrable [34, Section 473D].

Assumption 4.1. Throughout this section we assume that $h: \mathbb{R}^d \to \mathbb{R}$ is defined almost everywhere, $K: \mathbb{R}^d \to \mathbb{R}$ and that that the convolution product h*K exists almost everywhere.

4.2. Stability theorems. We now give several conditions on a pair of functions which imply that their convolution product is (locally) Lipschitz.

The first result appears in [34, 473D(d)], but the proof is included here for completeness.

Theorem 4.2. If $\|h\|_1 = a$ and K is b-Lipschitz, then h * K is ab-Lipschitz.

Proof. Let
$$g = h * K$$
. First we have, $g(u) - g(v) = \int_{\mathbb{R}^n} h(s) \left(K(u-s) - K(v-s) \right) ds$. Then, $|g(u) - g(v)| \leqslant \int_{\mathbb{R}^n} |h(s)| |K(u-s) - K(v-s)| ds \leqslant \int_{\mathbb{R}^n} |h(s)| b| u - v| ds \leqslant ab| u - v|$.

Let $B_{\alpha}(x)$ denote the closed ball of radius α centered at $x \in \mathbb{R}^d$, and let V_d denote the volume of the d-dimensional ball of radius 1.

Theorem 4.3. Let $x \in \mathbb{R}^d$ and let $\alpha > 0$. If $\|h\|_{\infty} \leqslant M$ on $B_{2\alpha}(x)$, K is b-Lipschitz and $supp(K) \subseteq B_{\alpha}(0)$, then h*K is $2Mb\alpha^dV_d$ -Lipschitz in $B_{\alpha}(x)$.

Proof. Let
$$g=h*K$$
. Let $u,v\in B_{\alpha}(x)$. As in the previous proof, $|g(u)-g(v)|\leqslant \int_{\mathbb{R}^n}|h(s)||K(u-s)-K(v-s)|ds\leqslant \int_{B_{\alpha}(u)\cup B_{\alpha}(v)}|h(s)|b|u-v|\,dx\leqslant 2Mb\alpha^dV_d|u-v|.$

Theorem 4.4. If $\|h\|_{\infty} \leq M$ and $\int |K(s+t) - K(s)| ds \leq b|t|$ for all $t \in \mathbb{R}^d$, then h * K is Mb-Lipschitz.

Proof. Let
$$g=h*K$$
. Again, $|g(u)-g(v)|\leqslant \int_{\mathbb{R}^n}|h(s)||K(u-s)-K(v-s)|ds\leqslant \int M|K(u-v+x)-K(x)|\,dx\leqslant Mb|u-v|.$

4.3. Application to kernels. We now apply the above theorems to smooth a function h, obtaining a Lipschitz function. That is, we will take K to be a *kernel*, a nonnegative integrable real-valued function on \mathbb{R}^n satisfying $\int K(x)dx = 1$, $\int xK(x)dx = 0$ and $\int x^2K(x)dx < \infty$. For example, we can choose K to be the *triangular kernel*, $K(x) = c \max(1 - \|x\|, 0)$, for appropriate normalization constant c. The most common choices

are the Gaussian kernel and the Epanechnikov kernel, which are described below. Notice that if K is a kernel, then so is $K_{\alpha}(x) = \frac{1}{\alpha^n} K(\frac{x}{\alpha})^2$. The parameter α is called the *bandwidth* and allows one to control the amount of smoothing.

4.3.1. The triangular kernel. Let $\alpha > 0$. Let V_d denote the volume of the n-dimensional ball of radius 1. For $A \subseteq \mathbb{R}^d$, let I_A denote the indicator function on A. That is, $I_A(x) = 1$ if $x \in A$ and 0 otherwise. The triangular kernel is given by

$$K_{\alpha}(x) = \frac{d+1}{\alpha^{d}V_{d}} \left(1 - \frac{|x|}{\alpha}\right) I_{B_{\alpha}(0)}.$$

Note that supp $(K_{\alpha}) = B_{\alpha}(0)$ and K_{α} is $\frac{d+1}{\alpha^{d+1}V_d}$ -Lipschitz. Applying Theorem 4.3, we have the following.

COROLLARY 4.5. Let $x \in \mathbb{R}^d$. If $\|h\|_{\infty} \leq M$ on $B_{2\alpha}(x)$ then $h * K_{\alpha}$ is $\frac{2M(d+1)}{\alpha}$ -Lipschitz in $B_{\alpha}(x)$.

Note that it follows that if the bound on h is global then so is the Lipschitz bound.

4.3.2. The Epanechnikov kernel. Let $\alpha > 0$. The Epanechnikov kernel is given by

$$K_{\alpha}(x) = \frac{d+2}{2\alpha^{d}V_{d}} \left(1 - \frac{|x|^{2}}{\alpha^{2}}\right) I_{B_{\alpha}(0)}.$$

Now supp $(K_{\alpha}) = B_{\alpha}(0)$ and K_{α} is $\frac{d+2}{\alpha^{d+1}V_d}$ -Lipschitz. Applying Theorem 4.3, we have the following.

COROLLARY 4.6. Let $x \in \mathbb{R}^d$. If $\|h\|_{\infty} \leq M$ on $B_{2\alpha}(x)$ then $h * K_{\alpha}$ is $\frac{2M(d+2)}{\alpha}$ -Lipschitz in $B_{\alpha}(x)$.

4.3.3. The Gaussian kernel. Let $\alpha > 0$. The Gaussian kernel is given by

$$K_{\alpha}(x) = \frac{1}{\alpha^d (2\pi)^{d/2}} e^{-|x|^2/2\alpha^2}.$$

Lemma 4.7. For the Gaussian kernel K_{α} , let $f(t)=\int |K_{\alpha}(s+t)-K_{\alpha}(s)|\,ds$. Then $f(t)\leqslant \frac{2}{\alpha\sqrt{2\pi}}|t|$ for all $t\in\mathbb{R}^d$.

²More generally, we can choose the bandwidth to be a symmetric positive definite matrix H and let $K_H(x) = \frac{1}{\sqrt{\det H}}K(H^{-1/2}x)$.

Proof. Change coordinates so that $s = -\frac{|t|}{2}e_1$ and $s + t = \frac{|t|}{2}e_1$. Then by symmetry

$$\begin{split} f(t) &= 2 \left[\int_{x_1 \geqslant -\frac{|t|}{2}} K_{\alpha}(x) \, dx - \int_{x_1 \geqslant \frac{|t|}{2}} K_{\alpha}(x) \, dx \right] \\ &= 4 \int_{0 \leqslant x_1 \leqslant \frac{|t|}{2}} K_{\alpha}(x) \, dx \\ &= \frac{4}{\alpha^d (2\pi)^{d/2}} \int_0^{\frac{|t|}{2}} e^{-x_1^2/2\alpha^2} \, dx_1 \int_{-\infty}^{\infty} e^{-x_2^2/2\alpha^2} \, dx_2 \cdots \int_{-\infty}^{\infty} e^{-x_d^2/2\alpha^2} \, dx_d \\ &= \frac{4}{\alpha \sqrt{2\pi}} \int_0^{\frac{|t|}{2}} e^{-x_1^2/2\alpha^2} \, dx_1 \end{split}$$

It follows that $f(t) \leqslant \frac{4}{\alpha \sqrt{2\pi}} \int_0^{|t|/2} dx_1 = \frac{2}{\alpha \sqrt{2\pi}} |t|.$

Thus by Theorem 4.4 we have the following

Corollary 4.8. If $\|h\|_{\infty} \leq M$ then $h * K_{\alpha}$ is $\frac{2M}{\alpha\sqrt{2\pi}}$ -Lipschitz.

In practice, the function h will rarely be essentially bounded, but this can be arranged by setting it to be 0 outside a closed ball centered at a specified configuration.

4.3.4. Sharpness of the Lipschitz constants. Assume that $K_\alpha:\mathbb{R}^d\to\mathbb{R}$ is symmetric in the first variable. Let $h:\mathbb{R}^d\to\mathbb{R}$ be defined by h(x)=1 if $x_1\geqslant 0$ and -1 otherwise. Let $g(t)=h*K_\alpha(te_1)-h*K_\alpha(-te_1)$. Then we calculate

$$\begin{aligned} h * K_{\alpha}(te_1) &= \int_{\mathbb{R}^d} h(te_1 - x) K_{\alpha}(x) \, dx \\ &= \int_{x_1 \leqslant t} K_{\alpha}(x) \, dx - \int_{x_1 > t} K_{\alpha}(x) \, dx \\ &= sign(t) \int_{-|t| \leqslant x_1 \leqslant |t|} K_{\alpha}(x) \, dx. \end{aligned}$$

So

$$g(t) = 2\, sign(t) \int_{-|t| \leqslant x_1 \leqslant |t|} K_\alpha(x) \ dx = 4\, sign(t) \int_{0 \leqslant x_1 \leqslant |t|} K_\alpha(x) \ dx.$$

Let K_{α} be the Gaussian kernel. Then

$$g(t) = \frac{4}{\alpha^{d}(2\pi)^{d/2}} \int_{0}^{t} e^{-x_{1}^{2}/2\alpha^{2}} dx_{1} \int_{-\infty}^{\infty} e^{-x_{2}^{2}/2\alpha^{2}} dx_{2} \cdots \int_{-\infty}^{\infty} e^{-x_{d}^{2}/2\alpha^{2}} dx_{d}$$

$$= \frac{4}{\alpha\sqrt{2\pi}} \int_{0}^{t} e^{-x_{1}^{2}/2\alpha^{2}} dx_{1}$$

It follows that g(t) converges to $\frac{4t}{\alpha\sqrt{2\pi}}$ as t approaches 0 by the first fundamental theorem of calculus. Hence, the Lipschitz constant given in Corollary 4.8 is optimal.

When d=1 and K_{α} is the triangular kernel, $g(t)=\frac{4\cdot 2}{\alpha V_1}\int_0^t (1-\frac{|x|}{\alpha})\,dx \to \frac{4t}{\alpha}$ as $t\to 0$. So the Lipschitz constant of $h*K_{\alpha}$ is at least $\frac{2}{\alpha}$. Hence, the Lipschitz constant given in Corollary 4.5 is optimal up to at most a factor of 2.

When d=1 and K_{α} is the Epanechnikov kernel, $g(t)=\frac{4\cdot 3}{2\alpha V_1}\int_0^t (1-\frac{x^2}{\alpha^2})\,dx \to \frac{3t}{\alpha}$ as $t\to 0$. So the Lipschitz constant of $h*K_{\alpha}$ is at least $\frac{3}{2\alpha}$. Hence, the Lipschitz constant given in Corollary 4.6 is optimal up to at most a factor of 4.

4.4. Stable Computations in Practice. Suppose that we can compute h(x) for values of x for which it is defined, we can sample from K, and that for a fixed $a \in \mathbb{R}^d$ we want to compute $g(a) = (h * K)(a) = \int_{\mathbb{R}^d} h(a-x)K(x)dx$. In practice, we will not be able to evaluate this integral analytically. We approximate g(a) as follows. Let V be a random variable with probability distribution given by the kernel K (one writes $V \sim K$). Let W be the random variable given by h(a-V). Then the expected value of W is given by $E[W] = \int_{\mathbb{R}^d} h(a-x)K(x)dx = g(a)$. We will approximate E[W] by drawing a sample e_1, \ldots, e_M where $e_i \sim K$ are independent. Then E[W] can be approximated by $\overline{W}_M = \frac{1}{M} \sum_{i=1}^M h(a-e_i)$. By the law of large numbers, $\overline{W}_M \to E[W]$, where the convergence may be taken to be in probability (the weak law) or almost surely (the strong law). This is the justification for the computations in Section 6. Let us record this result.

Theorem 4.9. Let $\alpha \in \mathbb{R}^d$ and $\varepsilon_1, \dots, \varepsilon_M$ be drawn independently from K. Then

$$\frac{1}{M}\sum_{i=1}^{M}h(\alpha-\epsilon_i)\to g(\alpha).$$

4.5. Stability of the choice of kernel. As should be clear, and as borne out by the experiments in Section 6, the value of (h * K)(a), for fixed h and a, will certainly depend on K. However, there is no fragility of output with respect to this choice, as shown by the following fact.

Theorem 4.10. Let $h: \mathbb{R}^d \to \mathbb{R}$ be an essentially bounded function. Then the map $K \to h*K$ is Lipschitz, from $L^1(\mathbb{R}^d)$ to $L^\infty(\mathbb{R}^d)$.

Proof. Let $\varphi:L^1(\mathbb{R}^d)\to L^\infty(\mathbb{R}^d)$ be given by $\varphi(K)=h*K.$ For $x\in\mathbb{R}^d$,

$$|\left[\varphi(K) - \varphi(K') \right](x)| \leqslant \int |h(x-t)| \, |K(t) - K'(t)| \, dt \leqslant \|h\|_{\infty} \|K - K'\|_{1}.$$

4.6. Choice of bandwidth. After choosing a family of kernels, such as the Gaussian kernels K_{α} described in Section 4.3.3, the most important choice in implementing the method described here is the choice of bandwidth α .

Choosing the amount of smoothing is a well-studied problem in nonparametric regression, where increasing the bandwidth decreases the estimation variance, but increases the squared bias. Both of these terms contribute to the error. A bandwidth which optimizes this trade-off may be estimated using cross-validation. A proper understanding of this problem in our situation requires analysis that goes beyond the scope of the present paper.

However, we offer some suggestions for the choice of bandwidth. First, it may be chosen to obtain a desired amount of smoothness of $h * K_{\alpha}$. For example, we may want $h * K_{\alpha}$ to be 1-Lipschitz. Second, it seems reasonable to choose the bandwidth to (at least) equal the level of estimated noise of the input data. One may combine these two to find the minimum bandwidth that satisfies both requirements.

- **4.7. Difficulties of working with persistence diagrams.** An underlying context for the work presented here is the difficulty of working with persistence diagrams.
- 4.7.1. Replacing persistence diagrams with features. Our approach centers on converting a persistence diagram to a real number. This may seem simplistic and somewhat ad hoc. However, all effective methods of combining persistence diagrams with serious statistical analysis and machine learning techniques rely on replacing a persistence diagram with a vector in some Hilbert space or Banach space. For simplicity, we restrict ourselves to the vector space \mathbb{R} , but our approach can be extended to more general vector spaces.
- 4.7.2. Repeated persistent homology computations. Applying Theorem 4.9 in Algorithm 1 we repeatedly compute persistence diagrams for similar filtrations. In the examples in Sections 1.2.1 and 1.2.2 we repeat M=1000 times. Note that we do not have convergence results at this time. For repeated persistent homology calculations it is important to have efficient software. In Section 1.2.2 we use Perseus [42] and in Section 1.2.1 we calculate persistent homology using a union-find data structure [32]. For other applications we would suggest Ripser [4].

In other applications, the computational cost may be considerable. We suggest the following. Consider the set of filtered complexes, L_i produced by $\{\alpha - \varepsilon_i\}$. If ε_i and ε_j are sufficiently close then L_i and L_j are isomorphic up to a change of filtration value and one doesn't need to repeat the persistence computation. More generally, one may obtain the persistent homology of similar filtered complexes using vineyard updates [27].

5. Application to persistent homology computations

Now let us apply the results of the previous section to persistent homology. Assume we have a persistent homology computation, \mathcal{C} , with input the real numbers α_1,\ldots,α_n . If our computation is defined for all $\alpha=(\alpha_1,\ldots,\alpha_n)\in\mathbb{R}^d$ and results in real values then we may proceed.

If not, we may reduce the more general situation to the one above as follows. We need h to be defined on all of \mathbb{R}^n so that convolutions with a Gaussian are well defined. Let \mathbb{O} be the set of outputs of this computation. Let $D\subseteq\mathbb{R}^n$ be the set of all inputs for

which $\mathbb C$ is defined. If $D \neq \mathbb R^n$ then add a state \emptyset to $\mathbb O$ and say that the computation sends all points in $\mathbb R^n - D$ to \emptyset . Thus we encode this computation as a function $H: \mathbb R^n \to \mathbb O$. Let $\mathfrak p$ be a real-valued function on $\mathbb O$ with $\mathfrak p(\emptyset) = \mathbb O$. Let $\mathfrak p = \mathbb P \cap \mathbb P$. We will need $\mathfrak p = \mathbb P \cap \mathbb P$. We will need $\mathfrak p = \mathbb P \cap \mathbb P$.

To make this less abstract, we show how the instabilities described in Sections 3.1 and 3.2 can be addressed by this method.

Example 5.1. Stable persistence located at a point. We return to Example 2.1, where we have a geometric line graph $\mathcal K$ with n vertices v_1,\ldots,v_n , and edges $e_i=(v_i,v_{i+1})$ for $i=1,\ldots n-1$. To produce a filtration of the type used in this example, we just need to know n function values. More precisely, our persistence computation takes as input a vector $\mathbf a=(a_1,\ldots,a_n)\in\mathbb R^n$, from which we obtain a piecewise linear function, F_a , on $\mathcal K$ determined by $F_a(v_i)=a_i$. Next we consider the corresponding abstract simplicial complex K and filtration f_a . Then we compute the persistence diagram $Dgm_0(f_a)$. This defines a function $H:\mathbb R^N\to \mathcal O$, where $H(a)=Dgm_0(f_a)$.

Now fix a specific vertex x in K. A given diagram $Dgm_0(f_\alpha)$ either contains a point u(x)=(b(x),d(x)) that represents a persistent connected component born at x in the filtration, or it does not. In the former case, we define $p_x(Dgm_0(f_\alpha))=d(x)-b(x)$, and in the latter we define $p(x)(Dgm_0(f_\alpha))=0$; that is, we map the diagram to the persistence of the connected component created by the addition of this specific vertex.

The discontinuity of the function $h_x = p_x \circ H : \mathbb{R}^n \to \mathbb{R}$ expresses the instability of localizing the persistence of a connected component. Referring to Figure 7, suppose that the vectors $\mathfrak a$ and $\mathfrak e$ produce the functions $\mathfrak f$ and $\mathfrak g$, respectively, and that the vertex $\mathfrak x$ is as marked in the figure. Then $h_x(\mathfrak a)$ is the persistence of $\mathfrak u$, while $h_x(\mathfrak e)$ is the persistence of $\mathfrak v'$. Corollaries 4.5 and 4.6 guarantee that smoothing h_x by a Triangular kernel and Epanechnikov kernel will result in a locally Lipschitz function.

To be able to convolve with the Gaussian kernel and apply Corollary 4.8 we need h_{σ} to be essentially bounded. We can arrange this by specifying that the domain D of $\mathcal C$ be compact and that p_{σ} be bounded. This requires that all of the persistence pairs in the output of $\mathcal C$ be finite. This can be arranged by truncating at some value M or by applying extended persistence [25]. The resulting h_x is Lipschitz.

The experiments in Section 6.1 show how this works in practice.

Example 5.2. Stable persistence located at an edge. We return to Example 2.2. In this case, K is the full complex on n vertices, and we start with n ordered points in the plane which lead to a piecewise-linear curve C. That is, $\mathcal C$ takes as input a vector $\mathfrak a \in \mathbb R^{2n}$ and places a vertex ν_i at $(\mathfrak a_{2i-1},\mathfrak a_{2i})$, thus creating a curve $C_\mathfrak a$. This leads to a filtration $f_\mathfrak a$ of K and finally we produce $Dgm_1(f_\mathfrak a) \in \mathcal O$. As before, $H(\mathfrak a) = Dgm_1(f_\mathfrak a)$ defines a function $H: \mathbb R^{2n} \to \mathcal O$.

If we fix a specific edge σ , we can proceed as in Example 5.1 by defining the function p_{σ} and thus $h_{\sigma} = p_{\sigma} \circ H$. For example, taking $\sigma = (A, D)$ in Figure 6 and letting α be the vector which led to that specific point configuration, we have $h_{\sigma}(\alpha)$ equal to the persistence of α . As above, α is (locally) Lipschitz.

Example 5.3. Stable persistence of generating cycles. Instead of tracking which j-simplex creates a persistent homology class, a persistent homology algorithm may record a j-cycle, γ , that represents the persistence class. In this case, we can define $p_{\gamma}: 0 \to \mathbb{R}$ to be d-b if γ represents a persistence pair [b,d) or otherwise 0. Let $h_{\gamma}=p_{\gamma}H$ and then $g_{\gamma}=h_{\gamma}*K_{\alpha}$ is (locally) Lipschitz.

Example 5.4. Stability in density-thresholding choice. Let Y be the point cloud on the bottom-left of Figure 8, which we recall was created from the point cloud on the top-left by adding three outlier points. Consider any de-noising process parametrized by some real numbers. For a specific example, let $k = (\delta, \varepsilon)$. For each $y \in Y$, let $C_{\delta}(y) = \{x \in Y \mid ||x - y|| \le \delta\}$. Then define

$$Y_\varepsilon^\delta = \{y \in Y \mid \frac{|C_\delta(y)|}{|Y|} \geqslant \varepsilon\}.$$

One then applies the Vietoris-Rips construction to obtain a filtered abstract simplicial complex from Y_{ε}^{δ} , and then computes $Dgm_1(Y_{\varepsilon}^{\delta})$. We may consider the input of our persistent homology computation ${\mathfrak C}$ to be $a_1,\dots,a_{2n},\delta,\varepsilon$: that is, the coordinates of the vertices and the parameter values. However, we may also take the coordinates to be fixed and only consider the parameters to be our input. Doing this, we obtain $H:{\mathbb R}^2\to {\mathbb O}$. In this case, define $p(D)=\max_{u\in D} pers(u)$ for any degree-one diagram D. Then the discontinuity of the function $h:{\mathbb R}^2\to{\mathbb R}$ given by

$$k = (\delta, \epsilon) \mapsto Dgm_1(Y_{\epsilon}^{\delta}) \mapsto p(Dgm_1(Y_{\epsilon}^{\delta}))$$

expresses the instability of the threshold-parameter choice referred to in Section 3.2. If k is chosen so that the three outlier points are remove, then h(k) will be the persistence of the most prominent point on the top-right of Figure 8. On the other hand, a very nearby choice of k might fail to remove these points, and we would get the persistence of the most prominent point on the bottom-right of Figure 8. As above, $g_{\sigma} = h_{\sigma} * K_{\alpha}$ is (locally) Lipschitz.

6. Experiments and Interpretations

This section more deeply investigates some of the examples above, both via a few proof-of-principle experiments with synthetic data and via some suggested interpretations of the results. All experiments were run in MATLAB, using TDATools for the persistent homology computations.

- **6.1. Experiments.** In this section we present three numerical experiments. The first two are instances of Example **5.1**: geometric line graphs in which we consider the statibility of persistence located at a point. The third is an instance of Example **5.2**: a piecewise linear curve in the plane in which we consider the stability of persistence at an edge.
- 6.1.1. First line-graph experiment. First we explore Example 5.1, where the input to a persistent homology computation is a choice of function-values on the vertices of a simplicial complex. Specifically, we consider a line graph $\mathcal K$ with vertices v_1,\ldots,v_7 , and the initial input choice $\alpha=(10,11,12.5,13,9.9,20,1)$. The left side of Figure 9 shows the graph of the PL-function F_α , and the persistence diagram $H(\alpha)$ is in the middle.

Note that the high-persistence dot (9.9,20) and the medium-persistence one (10,13) are created by the additions of ν_5 and ν_1 , respectively; that is, $h_5(\alpha)=10.01$ and $h_1(\alpha)=3$. These values are of course unstable to perturbations of α : for instance, if we switch the first and fifth entries of α , the reader can check that $h_5((9.9,11,12.5,13,10,20,1))=3$. Let K_{α} be a seven-dimensional Gaussian kernel with mean at the origin and bandwidth

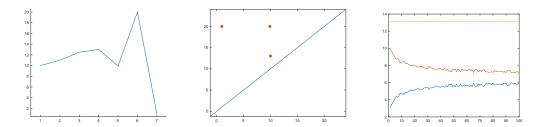


Figure 9. Input and results of first line graph experiment Left: a graph of the function F_{α} defined on a line graph with seven vertices. Middle: the persistence diagram $H(\alpha) = Dgm_0(f_{\alpha})$. We follow the extended persistence convention and pair the global min with the global max. Right: results of the experiment. Middle graph shows the value of $g_{5,\alpha}(\alpha)$ versus $10^3\alpha$; bottom graph shows $g_{1,\alpha}(\alpha)$ versus $10^3\alpha$; top graph shows their sum.

 α . For each $i=1,\ldots,7$, put $g_{i,\alpha}=h_i*K_\alpha$. The right side of Figure 9 shows graphs of the approximate values of $g_{5,\alpha}(\alpha)$ and $g_{1,\alpha}(\alpha)$, plotted against α , as well as a graph of their sum. There were 100 evenly spaced values of α used, ranging from $\alpha=0.001$ to $\alpha=0.1$. To make these graphs, we followed the approximation procedure suggested by Theorem 4.9. For each fixed α , we took N=1000 independent draws $\varepsilon_1,\ldots\varepsilon_{1000}$ from K_α , and computed

$$g_5(\alpha) \approx \frac{1}{1000} \sum_{i=1}^{1000} h_5(\alpha + \epsilon_i),$$

with an identical procedure for $g_1(a)$.

6.1.2. Second line-graph experiment. Again we explore Example 5.1, this time with the input a=(5,1.1,1,1.05,15) to a persistent homology computation that builds a filtration on a line graph with five vertices $v_1,\ldots v_5$. The function F_a , whose graph is on the left of Figure 10, has a global min at v_3 . From the diagram in the middle, we see $h_3(a)=15-1=14$. Note that $h_i(a)=0$ for $i\neq 3$, since only one component is created during the entire filtration. On the right, we see convolved values of these functions, with notation and computation procedure exactly as in 6.1.1 above.

6.1.3. Distance-to-a-curve experiment. Next we reconsider Example 5.2. Let C be the PL-curve with nine vertices on the left side of Figure 11. In our language, $C = C_a$, where the input vector a specifies the coordinates of the nine vertices: $v_1 = (0,0.1), v_2 = (1,1), v_3 = (2,0.12), v_4 = (7,5), v_5 = (12,0), v_6 = (7,-5), v_7 = (2,-0.12), v_8 = (1,-1), v_9 = (0,-0.1)$. Following the vocabulary of Example 2.2, this curve placement leads to a order-preserving function f_a on the abstract full complex K on nine vertices.

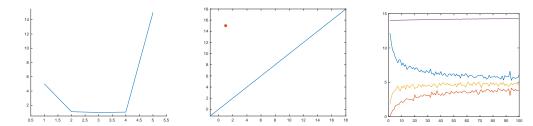


FIGURE 10. Input and results of second line graph experiment Left: a graph of the function F_a defined on a line graph with seven vertices. Middle: the persistence diagram $H(a) = Dgm_0(f_a)$. We follow the extended persistence convention and pair the global min with the global max. Right: results of experiment. Moving from bottom to top, the values of $g_{2,\alpha}(a)$, $g_{4,\alpha}(a)$, $g_{3,\alpha}(a)$, and their sum, all plotted against $10^3\alpha$.

Its degree-one persistence diagram $H(a) = Dgm_1(f_a)$, in the middle of the same figure, has only two off-diagonal points. The first, at (0.2, 10), is created by the positive edge between v_1 and v_9 , while the second, at (0.23, 2), comes from the edge between v_3 and v_7 . Thus we have $h_{1,9}(a) = 9.8$ and $h_{3,7}(a) = 1.77$. As usual, these values are highly unstable to small perturbations in the vertex positions. In very similar fashion to the

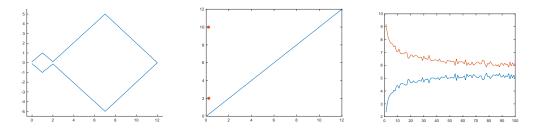


Figure 11. Input and results of distance-to-curve experiment. Left: the PL plane curve C_{α} whose nine vertices are defined in the text. Middle: the persistence diagram $H(\alpha) = Dgm_1(f_{\alpha})$. Right: Results of experiment. Top graph shows the value of $g_{1,9,\alpha}(\alpha)$ versus $10^4\alpha$, bottom graph shows $g_{3,7,\alpha}(\alpha)$ versus $10^4\alpha$.

last experiment, we then computed approximate values at α for the convolved functions $g_{1,9,\alpha}=f_{1,9}*K_{\alpha}$ and $g_{3,7,\alpha}=f_{3,7}*K_{\alpha}$ where K_{α} was a nine-dimensional Gaussian kernel with bandwidth α . This time we used 100 evenly spaced values of α , going from 0.0001 to 0.01. The results appear on the right side of Figure 11.

- **6.2. Interpretations.** We now offer some possible interpretations one can draw from these results, and also suggest some potential uses of this technique in practice.
- 6.2.1. Locating a point in the domain. Let u = (9.9, 20) be one of the high-persistence points in the diagram for our first line-graph experiment. It is accurate to say that u was created, for this specific persistent homology computation, by the addition of v_5 to the filtration. However, it is also a potentially misleading thing to say.

We propose that the difference between the persistence of u and the values of the convolutions $g_{5,\alpha}(a)$ might be seen as an indicator for how confidently one should locate u at v_5 . The graphs on the right side of Figure 9 tell us that this confidence should be low. On the other hand, the other high-persistence point w = (1, 20) is created by the addition of v_7 . It turns out that $g_{7,\alpha}(a)$ remains very close to 19 for all α within a reasonable range.

- 6.2.2. Spreading out a point in the domain. Alternatively, one might choose to give u a more fuzzy location. A reasonable idea would be to spread out its location between vertices v_5 and v_1 , since v_1 is responsible for creating the same component in a nearby filtration. The graphs in figure 9 bear this out: note that the sum of the two convolution values $g_{5,\alpha}(\alpha) + g_{1,\alpha}(\alpha)$ is always very close to the sum of the persistences of the components created by v_1 and v_5 . Similarly, in the second line-graph experiment, it would be reasonable to smear the location of the only point throughout the immediate neighborhood of v_3 .
- 6.2.3. Convolved values as features. One could also use the values of g_i or $g_{i,j}$ as features in a machine-learning scheme. That is, the vector $(g_{1,\alpha}(a), \ldots, g_{7,\alpha}(a))$ could be used as a summary feature of both the filtration created by α and the noise model K_{α} . The stabilities offered by Corollary 4.8 and Theorem 4.10 make this an appealing option.

7. Future work

The work presented here opens many interesting questions and possible directions for future research.

Machine learning. Persistence diagrams have been used to produce features for machine-learning and statistical methods. This paper takes a first step towards the extraction of stable features that describe much of the other information produced during a persistent homology computation. For example, one could apply the ideas presented here to construct topological features for machine learning that are not only based on critical values but also on the locations of critical points.

Visualization. It would be nice to use the ideas of this paper to build visualization tools. For example, one might want to compute a persistence diagram, click on a point, and have the possible location candidates shown on the domain, perhaps with some sort of heat map of likelihood.

Convergence results. It would be good to have theoretical results on the rate of convergence of Algorithm 1 as the number of repetitions M increases.

Samples with increasing numbers of points. In all of the examples we consider, the number of sampled points is fixed. It would be nice to expand this framework to allow for increasing sample sizes, allowing asymptotic results to be considered.

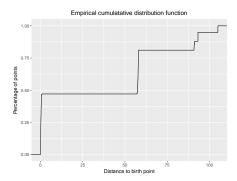


FIGURE 12. The empirical cumulative distribution function of the distance between the location of the generator of the 28th longest bar in one of the iterations of Algorithm 1 to the location of the generator of the 28th longest bar in the observed data.

A continuous theory. Finally, we also hope to enrich the theory whose development has started here. It would be nice to work instead in the category of metric spaces, to define some versions of the functions h_x and g_x , and to prove Lipschitz-continuity of the latter. We believe that the direction suggested by Example 2.6 may be the right one.

Features in other vector spaces. The present paper has only considered the stabilization of real-valued functions. However, one could consider functions in \mathbb{R}^n or more generally into Banach spaces.

Acknowledgments

The authors would like to thank Justin Curry, Francis Motta, Chris Tralie, and Ulrich Bauer for helpful conversations. The first author would like to thank the University of Florida for hosting him during the initial research phase. We would also like to thank the referees for their careful reading of the paper and their helpful suggestions.

Appendix A. Brain arteries

We can obtain the estimate of h * K for the observed data in Section 1.2.1 for a ball of any radius by considering the distance from the location of the generator of the 28th longest bar in one of the iterations of Algorithm 1 to the location of the generator of the 28th longest bar in the observed data, and computing the empirical cumulative distribution function of this function. See Figure 12. In this figure we see that there is a competitor for the location of the generator of the 28th longest bar. So in fact, the situation with this real data is quite similar to that in the elementary synthetic examples in Sections 6.1.1 and 6.1.3.

References

- [1] Henry Adams, Atanas Atanasov, and Gunnar Carlsson. Nudged elastic band in topological data analysis. *Topol. Methods Nonlinear Anal.*, 45(1):247–272, 2015.
- [2] Henry Adams, Tegan Emerson, Michael Kirby, Rachel Neville, Chris Peterson, Patrick Shipman, Sofya Chepushtanova, Eric Hanson, Francis Motta, and Lori Ziegelmeier. Persistence images: A stable vector representation of persistent homology. *Journal of Machine Learning Research*, 18(8):1–35, 2017.
- [3] Mahmuda Ahmed, Brittany Terese Fasy, and Carola Wenk. Local persistent homology based distance between maps. In *SIGSPATIAL*. ACM, Nov. 2014.
- [4] Ulrich Bauer. Ripser: a lean C++ code for the computation of vietoris—rips persistence barcodes. Software available at https://github.com/Ripser/ripser, 2017.
- [5] P. Bendich, S. P. Chin, J. Clark, J. Desena, J. Harer, E. Munch, A. Newman, D. Porter, D. Rouse, N. Strawn, and A. Watkins. Topological and statistical behavior classifiers for tracking applications. *IEEE Transactions on Aerospace and Electronic Systems*, 52(6):2644–2661, December 2016.
- [6] Paul Bendich, J. S. Marron, Ezra Miller, Alex Pieloch, and Sean Skwerer. Persistent homology analysis of brain artery trees. *Ann. Appl. Stat.*, 10(1):198–218, 03 2016.
- [7] Andrew J. Blumberg, Itamar Gal, Michael A. Mandell, and Matthew Pancia. Robust statistics, hypothesis testing, and confidence intervals for persistent homology on metric measure spaces. *Found. Comput. Math.*, 14(4):745–789, 2014.
- [8] Peter Bubenik. Statistical topological data analysis using persistence landscapes. *Journal of Machine Learning Research*, 16:77–102, 2015.
- [9] Peter Bubenik, Gunnar Carlsson, Peter T. Kim, and Zhi-Ming Luo. Statistical topology via Morse theory persistence and nonparametric estimation. In *Algebraic Methods in Statistics and Probability II*, volume 516 of *Contemp. Math.*, pages 75–92. Amer. Math. Soc., Providence, RI, 2010.
- [10] Gunnar Carlsson. Topology and data. Bulletin of the American Mathematical Society, 46(2):255–308, January 2009.
- [11] Gunnar Carlsson, Tigran Ishkhanov, Vin de Silva, and Afra Zomorodian. On the local behavior of spaces of natural images. *Int. J. Comput. Vision*, 76(1):1–12, 2008.
- [12] Gunnar Carlsson, Tigran Ishkhanov, Vin de Silva, and Afra Zomorodian. On the local behavior of spaces of natural images. *International Journal of Computer Vision*, 76(1):1–12, 2008.
- [13] M. Carriere, S. Y. Oudot, and M. Ovsjanikov. Stable topological signatures for points on 3d shapes. In *Proc. Sympos. on Geometry Processing*, 2015.
- [14] Frédéric Chazal, David Cohen-Steiner, Marc Glisse, Leonidas J. Guibas, and Steve Y. Oudot. Proximity of persistence modules and their diagrams. In *Proceedings of the 25th annual symposium on Computational geometry*, SCG '09, pages 237–246, New York, NY, USA, 2009. ACM.
- [15] Frédéric Chazal, David Cohen-Steiner, and Quentin Mérigot. Geometric Inference for Measures based on Distance Functions. *Foundations of Computational Mathematics*, 11(6):733–751, 2011. RR-6930 RR-6930.
- [16] Frédéric Chazal, Vin de Silva, Marc Glisse, and Steve Oudot. *The structure and stability of persistence modules*. SpringerBriefs in Mathematics. Springer, [Cham], 2016.
- [17] Frédéric Chazal, Brittany T. Fasy, Fabrizio Lecci, Bertrand Michel, Alessandro Rinaldo, and Larry Wasserman. Robust topological inference: Distance to a measure and kernel distance. 12 2014, 1412.7197.
- [18] Frédéric Chazal, Brittany Terese Fasy, Fabrizio Lecci, Bertrand Michel, Alessandro Rinaldo, and Larry Wasserman. Subsampling methods for persistent homology. In *Proceedings of the 32nd International Conference on Machine Learning, Lille, France*, volume 37. JMLR: W&CP, 2015.
- [19] Frédéric Chazal, Brittany Terese Fasy, Fabrizio Lecci, Alessandro Rinaldo, Aarti Singh, and Larry Wasserman. On the bootstrap for persistence diagrams and landscapes. *Modeling and Analysis of Information Systems*, 20(6):96–105, 2014.
- [20] Frédéric Chazal, Brittany Terese Fasy, Fabrizio Lecci, Alessandro Rinaldo, and Larry Wasserman. Stochastic convergence of persistence landscapes and silhouettes. *J. Comput. Geom.*, 6(2):140–161, 2015.
- [21] Y.-C. Chen, D. Wang, A. Rinaldo, and L. Wasserman. Statistical analysis of persistence intensity functions. *ArXiv e-prints*, 2015, 1510.02502.

- [22] Yen-Chi Chen, Christopher R. Genovese, and Larry Wasserman. Statistical inference using the morse-smale complex. 06 2015.
- [23] Moo K. Chung, Peter Bubenik, and Peter T. Kim. Persistence diagrams in cortical surface data. In *Information Processing in Medical Imaging (IPMI)* 2009, volume 5636 of *Lecture Notes in Computer Science*, pages 386–397, 2009.
- [24] David Cohen-Steiner, Herbert Edelsbrunner, and John Harer. Stability of persistence diagrams. *Discrete Comput. Geom.*, 37(1):103–120, January 2007.
- [25] David Cohen-Steiner, Herbert Edelsbrunner, and John Harer. Extending persistence using Poincaré and Lefschetz duality. *Found. Comput. Math.*, 9(1):79–103, 2009.
- [26] David Cohen-Steiner, Herbert Edelsbrunner, John Harer, and Yuriy Mileyko. Lipschitz functions have lp-stable persistence. *Found. Comput. Math.*, 10(2):127–139, February 2010.
- [27] David Cohen-Steiner, Herbert Edelsbrunner, and Dmitriy Morozov. Vines and vineyards by updating persistence in linear time. In *Computational geometry (SCG'06)*, pages 119–126. ACM, New York, 2006.
- [28] Vin de Silva and Gunnar Carlsson. Topological estimation using witness complexes. *Eurographics Symposium on Point-Based Graphics*, 2004.
- [29] Tamal K. Dey and Rephael Wenger. Stability of critical points with interval persistence. *Discrete Comput. Geom.*, 38(3):479–512, 2007.
- [30] Tamal Krishna Dey, Fengtao Fan, and Yusu Wang. Graph induced complex on point data. In *Proceedings of the Twenty-ninth Annual Symposium on Computational Geometry*, SoCG '13, pages 107–116, New York, NY, USA, 2013. ACM.
- [31] Herbert Edelsbrunner and John Harer. *Computational Topology: An Introduction*. American Mathematical Society, 2010.
- [32] Herbert Edelsbrunner and John L. Harer. *Computational topology*. American Mathematical Society, Providence, RI, 2010.
- [33] Brittany Terese Fasy, Fabrizio Lecci, Alessandro Rinaldo, Larry Wasserman, Sivaraman Balakrishnan, and Aarti Singh. Confidence sets for persistence diagrams. *Ann. Statist.*, 42(6):2301–2339, 2014.
- [34] D.H. Fremlin. Measure Theory, Volume 4. Torres Fremlin, 2000.
- [35] P. Frozini and B. Landi. Size theory as a topological tool for computer vision. *Pattern Recognition and Image Analysis*, **9**:596–603, 1999.
- [36] Marcio Gameiro, Yasuaki Hiraoka, Shunsuke Izumi, Miroslav Kramar, Konstantin Mischaikow, and Vidit Nanda. A topological measurement of protein compressibility. *Jpn. J. Ind. Appl. Math.*, 32(1):1–17, 2015.
- [37] Robert Ghrist. Barcodes: the persistent topology of data. Bull. Amer. Math. Soc. (N.S.), 45(1):61–75, 2008.
- [38] Violeta Kovacev-Nikolic, Peter Bubenik, Dragan Nikolić, and Giseon Heo. Using persistent homology and dynamical distances to analyze protein binding. *Stat. Appl. Genet. Mol. Biol.*, 15(1):19–38, 2016.
- [39] Chunyuan Li, M. Ovsjanikov, and F. Chazal. Persistence-based structural recognition. In *Computer Vision and Pattern Recognition (CVPR)*, 2014 IEEE Conference on, pages 2003–2010, June 2014.
- [40] Elizabeth Munch, Katharine Turner, Paul Bendich, Sayan Mukherjee, Jonathan Mattingly, John Harer, et al. Probabilistic fréchet means for time varying persistence diagrams. *Electronic Journal of Statistics*, 9:1173–1204, 2015.
- [41] James R. Munkres. Elements of Algebraic Topology. Addison Wesley, 1993.
- [42] Vidit Nanda. Perseus: the persistent homology software. Software available at http://www.math.rutgers.edu/~vidit/perseus/index.html, 2013.
- [43] P. Niyogi, S. Smale, and S. Weinberger. A topological view of unsupervised learning from noisy data. *SIAM J. Comput.*, 40(3):646–663, 2011.
- [44] Steve Y. Oudot. *Persistence theory: from quiver representations to data analysis*, volume 209 of *Mathematical Surveys and Monographs*. American Mathematical Society, Providence, RI, 2015.
- [45] Jan Reininghaus, Stefan Huber, Ulrich Bauer, and Roland Kwitt. A stable multi-scale kernel for topological machine learning. *Proc. CVPR*, pages 4741–4748, 2015.
- [46] B. W. Silverman. *Density estimation for statistics and data analysis*. Monographs on Statistics and Applied Probability. Chapman & Hall, London, 1986.
- [47] Abraham Smith, Paul Bendich, John Harer, and Jay Hineman. Supervised learning of labeled point-cloud differences via cover-tree entropy reduction. 02 2017.

- [48] M. P. Wand and M. C. Jones. *Kernel smoothing*, volume 60 of *Monographs on Statistics and Applied Probability*. Chapman and Hall, Ltd., London, 1995.
- [49] Shmuel Weinberger. The complexity of some topological inference problems. *Found. Comput. Math.*, 14(6):1277–1285, 2014.
- [50] Afra Zomorodian and Gunnar Carlsson. Localized homology. Comput. Geom., 41(3):126–148, 2008.

DEPARTMENT OF MATHEMATICS, DUKE UNIVERSITY, AND GEOMETRIC DATA ANALYTICS, INC.

E-mail address: bendich@math.duke.edu

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF FLORIDA

E-mail address: peter.bubenik@ufl.edu

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF FLORIDA

E-mail address: wagnera@ufl.edu