

Finding the Homology of Submanifolds with High Confidence from Random Samples*

Partha Niyogi,¹ Stephen Smale,² and Shmuel Weinberger³

¹Departments of Computer Science and Statistics, University of Chicago,
Chicago, IL 60637, USA
niyogi@cs.uchicago.edu

²Toyota Technological Institute, University Press Building,
Chicago, IL 60637, USA
smale@tti-c.org

³Department of Mathematics, University of Chicago,
Chicago, IL 60637, USA
shmuel@math.uchicago.edu

Abstract. Recently there has been a lot of interest in geometrically motivated approaches to data analysis in high-dimensional spaces. We consider the case where data is drawn from sampling a probability distribution that has support on or near a submanifold of Euclidean space. We show how to “learn” the homology of the submanifold with high confidence. We discuss an algorithm to do this and provide learning-theoretic complexity bounds. Our bounds are obtained in terms of a condition number that limits the curvature and nearness to self-intersection of the submanifold. We are also able to treat the situation where the data is “noisy” and lies near rather than on the submanifold in question.

1. Introduction

In recent years there has been considerable interest in the possibility of analyzing and processing data in high-dimensional spaces. Following the intuition that naturally occurring data may be generated by structured systems with possibly much fewer degrees of freedom than the ambient dimension would suggest, various researchers (see [16], [17], [3], [10], and [20]) have considered the case when the data lives on or close to a

* The main results of this paper were first presented at a conference in honor of John Franks and Clark Robinson at Northwestern University in April 2003. These results were formally written as Technical Report No. TR-2004-08, Department of Computer Science, University of Chicago.

submanifold of the ambient space. One hopes then to estimate geometrical and topological properties of the submanifold from random points (“scattered data”) lying on this unknown submanifold. These questions belong to a class of problems that have come to be known as *manifold learning*.

In this paper we consider the particular question of identifying the homology of the submanifold from random samples. The homology of the submanifold (see [15] for definitions) are natural topological invariants that provide a good characterization of many aspects of it. For example, the dimensions of the homology groups, the Betti numbers (b_0, b_1, \dots) , have natural interpretations. b_0 , the dimension of the zeroth homology group is the number of connected components of the submanifold. In data analysis situations, the number of clusters of the data may sometimes be understood in terms of the number of components of an underlying manifold (or other geometric object). If the dimension of the submanifold is d , then one sees that $b_j = 0$ for all $j > d$. Thus the the largest non-trivial homology gives us the dimension of the submanifold. If the submanifold is two-dimensional, then b_0 and b_1 are related to the number of connected components and number of holes, respectively, of the submanifold.

We show that it is possible to identify the homology from random samples and discuss an algorithm to do this. There are a few aspects of the developments in this paper that are worth emphasizing. First, we provide sample complexity estimates on the number of examples that are needed to identify the homology with high confidence. Our results are in the style of learning-theoretic treatments (for example, the *Probably Approximately Correct* framework [18]) where unknown objects (typically functions in learning theory) are “learned” from random samples and confidence estimates are provided. Second, we treat the situation where data might be drawn from a distribution that is concentrated *around* the manifold rather than precisely on it. Under specific models of noise, we show that our algorithm can work even with noisy data. In all cases, estimates are provided in terms of a condition number that limits the curvature and nearness to self-intersection of the submanifold.

Our results may also be of interest to researchers in computational geometry and topology who have considered the question of computing homology from simplicial complexes in the past (see [14] and [8] for details and further references). A number of researchers in these computational geometry and topology fields have considered the problem of manifold reconstruction from point cloud data. Such work has typically focused on the case of surfaces in \mathbb{R}^3 and examples include algorithms associated with the frameworks of alpha shapes [11], CRUST [1] and its variants, and COCONE [2] and its generalizations. CRUST and COCONE provably recover a simplicial 2-manifold that is homeomorphic to the surface. In [6] (written after the results of our current paper were declared), it was shown how to extend these ideas to the general setting of a k -manifold embedded in \mathbb{R}^N . In much of this work the medial axis plays a central role in characterizing the conditioning of the manifold (see our later remarks in Section 2). It is also worth noting that none of the works mentioned above considers the probabilistic setting where examples are drawn at random—so no high confidence guarantees are provided. The theorems in [1], [2], and [6] are analogous to our Proposition 3.1. No version of our main theorem (Theorem 3.1) exists in the literature. Finally, it is also worth noting that there is a body of work on persistence homology [20], [7] that seeks alternative topological characterizations of the manifold and its homology. See the discussion after Proposition 3.1.

In conclusion, we hope that researchers in graphics, pattern recognition, solid modeling, molecular biology, finance, and other areas where large amounts of high-dimensional data are available may find some use for the topological perspective on data analysis embodied in the algorithms and analyses of this paper.

2. Preliminaries

Consider a compact Riemannian submanifold \mathcal{M} of a Euclidean space \mathbb{R}^N . Sample the manifold according to a uniform probability measure on it. Thus points $x_1, \dots, x_n \in \mathcal{M}$ are generated. This set of points $\bar{x} = \{x_1, \dots, x_n\}$ is the data set on the basis of which homology groups will be calculated. In later sections we consider the case when the data is drawn from a probability measure with support close to the manifold.

Throughout our discussion, we associate to \mathcal{M} a condition number $(1/\tau)$ where τ is defined as the largest number having the property: The open normal bundle about \mathcal{M} of radius r is embedded in \mathbb{R}^N for every $r < \tau$. Its image Tub_τ is a tubular neighborhood of \mathcal{M} with its canonical projection map

$$\pi_0: \text{Tub}_\tau \rightarrow \mathcal{M}.$$

Note that τ encodes both local curvature considerations as well as global ones: If \mathcal{M} is a union of several components, then τ bounds their separation. For example, if \mathcal{M} is a sphere, then τ is equal to its radius. If \mathcal{M} is an annulus, then τ is the separation of its components. In Section 6 we relate the condition number $1/\tau$ to classical notions of curvature in differential geometry via the second fundamental form.

Finally, it is also useful to relate τ to the notions of medial axis and local feature size that have been developed in the computational geometry community. Given \mathcal{M} , one may define the set

$$G = \{x \in \mathbb{R}^N \text{ such that } \exists \text{ distinct } p, q \in \mathcal{M} \text{ where } d(x, \mathcal{M}) = \|x - p\| = \|x - q\|\},$$

where $d(x, \mathcal{M}) = \inf_{y \in \mathcal{M}} \|x - y\|$ is the distance of x to \mathcal{M} . The closure of G is called the medial axis and for any point $p \in \mathcal{M}$ the local feature size $\sigma(p)$ is the distance of p to the medial axis. Then it is easy to check that

$$\tau = \inf_{p \in \mathcal{M}} \sigma(p).$$

3. An Outline of Our Main Results

Ultimately we wish to compute the homology of the manifold $\mathcal{M} \subset \mathbb{R}^N$ from the randomly sampled datapoints $\bar{x} = \{x_1, \dots, x_n\} \subset \mathcal{M}$. We first begin by considering Euclidean balls (in the ambient space \mathbb{R}^N) of radius ε and center x_i . We denote these balls as $B_\varepsilon(x_i)$. We can now define the open set $U \subset \mathbb{R}^N$ given by

$$U = \bigcup_{x \in \bar{x}} B_\varepsilon(x).$$

Our first proposition states that if $\bar{x} = \{x_1, \dots, x_n\}$ is $\varepsilon/2$ dense in \mathcal{M} , then \mathcal{M} is a deformation retract of U .

Proposition 3.1. *Let \bar{x} be any finite collection of points $x_1, \dots, x_n \in \mathbb{R}^N$ such that it is $(\varepsilon/2)$ dense in \mathcal{M} , i.e., for every $p \in \mathcal{M}$, there exists an $x \in \bar{x}$ such that $\|p-x\|_{\mathbb{R}^N} < \varepsilon/2$. Then for any $\varepsilon < \sqrt{\frac{3}{5}}\tau$, we have that U deformation retracts to \mathcal{M} . Therefore the homology of U equals the homology of \mathcal{M} .*

We prove this proposition in Section 4. Subsequent to our work, the authors of [7] presented a different type of calculation of the homology of \mathcal{M} based on their homology approximation theorem together with the method of computing persistent homology (e.g., [20]). Their method does not give the homotopy type of \mathcal{M} . On the other hand, it does apply to a class of metric spaces more general than well-conditioned manifolds. A related approach appears in [5].

In the case under consideration here, the points x_1, \dots, x_n are sampled in i.i.d. fashion from the uniform probability distribution on \mathcal{M} . By probabilistic considerations, we will then prove (in Section 5)

Proposition 3.2. *Let \bar{x} be drawn by sampling \mathcal{M} in i.i.d. fashion according to the uniform probability measure on \mathcal{M} . Then with probability greater than $1 - \delta$, we have that \bar{x} is $(\varepsilon/2)$ -dense ($\varepsilon < \tau/2$) in \mathcal{M} provided*

$$|\bar{x}| > \beta_1 \left(\log(\beta_2) + \log\left(\frac{1}{\delta}\right) \right),$$

where

$$\beta_1 = \frac{\text{vol}(\mathcal{M})}{(\cos^k(\theta_1)) \text{vol}(B_{\varepsilon/4}^k)} \quad \text{and} \quad \beta_2 = \frac{\text{vol}(\mathcal{M})}{(\cos^k(\theta_2)) \text{vol}(B_{\varepsilon/8}^k)}.$$

Here k is the dimension of the manifold \mathcal{M} and $\text{vol}(B_\varepsilon^k)$ denotes the k -dimensional volume of the standard k -dimensional ball of radius ε . Finally, $\theta_1 = \arcsin(\varepsilon/8\tau)$ and $\theta_2 = \arcsin(\varepsilon/16\tau)$.

Putting these two propositions together, we see that we are able to provide a finite sample estimate for how many times we need to sample \mathcal{M} so that we are guaranteed with high confidence that the homology of the random set U equals the homology of \mathcal{M} . Thus our main theorem is

Theorem 3.1. *Let \mathcal{M} be a compact submanifold of \mathbb{R}^N with condition number τ . Let $\bar{x} = \{x_1, \dots, x_n\}$ be a set of n points drawn in i.i.d. fashion according to the uniform probability measure on \mathcal{M} . Let $0 < \varepsilon < \tau/2$. Let $U = \bigcup_{x \in \bar{x}} B_\varepsilon(x)$ be a correspondingly random open subset of \mathbb{R}^N . Then for all*

$$n > \beta_1 \left(\log(\beta_2) + \log\left(\frac{1}{\delta}\right) \right),$$

the homology of U equals the homology of \mathcal{M} with high confidence (probability $> 1 - \delta$).

Remark. Note that no version of our main theorem exists in the literature so far. However, versions of our Proposition 3.1 do exist. We have characterized Proposition 3.1 in terms of τ but one may obtain an alternate characterization in terms of the medial axis and the local feature size. In fact, if one considers the union of balls centered at the data points given by $U = \bigcup_{x \in \bar{x}} B_{\varepsilon_x}(x)$ where $\varepsilon_x = r\sigma(x)$, then it is possible to show that the homology of U coincides with that of \mathcal{M} if \bar{x} is $(\varepsilon_x/2)$ -dense in \mathcal{M} and for all $r < 0.21$. For the case of surfaces in \mathbb{R}^3 , a similar result is obtained by Amenta et al. [2] for $r < 0.06$. The set \bar{x} is said to be $(\varepsilon_x/2)$ -dense if for every $p \in \mathcal{M}$ there exists some $x \in \bar{x}$ such that $\|p - x\| < \varepsilon_x/2$. We will prove this in a later paper. It is not obvious, however, how to obtain a version of our main theorem in terms of the local feature size. Finally, we recall the recent results of [7] that we have already alluded to.

3.1. Computing the Homology of U

One now needs to consider algorithms to compute the homology of U . Noting that the $B_\varepsilon(x_i)$'s form a cover of U , one can construct the *nerve* of the cover. The nerve is an abstract simplicial complex constructed as follows: One puts in a k -simplex for every $(k + 1)$ -tuple of intersecting elements of the cover. The Nerve Lemma (see [4]) applies in our case, as balls are convex, to show that the homology of U is the same as the homology of this complex. The algorithm consists of the following components:

1. Given an ε , and a set of points $\bar{x} = \{x_1, \dots, x_n\}$ in \mathbb{R}^N , each j -simplex is given by a subset of the n points that have non-zero intersection. Thus we may define L_j to be the collection of all j -simplices. Each simplex $\sigma \in L_j$ is associated with a set of $j + 1$ points $(p_0(\sigma), \dots, p_j(\sigma) \in \bar{x})$ such that

$$\bigcap_{i=0}^j B_\varepsilon(p_i(\sigma)) \neq \emptyset.$$

An orientation for the simplex is chosen by picking an ordering and we denote the oriented simplex by $|p_0(\sigma), \dots, p_j(\sigma)|$.

2. A very crude upper bound on the size of L_j (denoted by $|L_j|$) is given by $\binom{n}{j+1}$. However, it is clear that if two points x_m and x_l are more than 2ε apart, they cannot be associated to a simplex. Therefore, there is a locality condition that the $p_i(\sigma)$'s must obey which results in $|L_j|$ being much smaller than this crude number. The simplicial complex $K_j = \bigcup_{i=0}^j L_j$ together with face relations. The simplicial complex corresponding to the nerve of U is $K = K_N$.
3. A basic subroutine for computing the simplicial complex (steps 1 and 2 above) involves the decision problem: for any set of j points, determine whether balls of radius ε around each of these points have non-empty intersection. This problem is related to the smallest ball problem defined as follows: Given a set of j points, find the ball with the smallest radius enclosing all these points. One can check that $\bigcap_{i=1}^j B_\varepsilon(p_i) \neq \emptyset$ if and only if this smallest radius $< \varepsilon$. Fast algorithms for the smallest ball problem exist. See [12] for theoretical discussion and [14] for downloadable algorithms from the web.

4. We work in the field of coefficients \mathbb{R} . Then a j -chain is a function $c: L_j \rightarrow \mathbb{R}$ and can be written as a formal sum

$$c = \sum_{\sigma \in L_j} c(\sigma)\sigma.$$

By adding j -chains componentwise, one gets the vector space of j -chains denoted by C_j .

5. The boundary operator ∂_j is a linear operator from C_j to C_{j-1} defined as follows. For each (oriented) simplex $\sigma \in L_j$,

$$\partial_j \sigma = \sum_{i=0}^j (-1)^i \sigma_i,$$

where σ_i is a $j-1$ face of σ (facing point $p_i(\sigma)$) and the orientation of σ_i is given by $|p_0, \dots, p_{i-1}, p_{i+1}, \dots, p_j|$. Now ∂_j is defined on j chains by additivity as

$$\partial_j \left(\sum_{\sigma \in L_j} c(\sigma)\sigma \right) = \sum_{\sigma \in L_j} c(\sigma)\partial_j \sigma.$$

Thus, ∂_j can be represented as an $n_{j-1} \times n_j$ matrix where $n_{j-1} = |L_{j-1}|$ and $n_j = |L_j|$, respectively. The matrix is usually sparse in our setting.

6. This defines the chain complex

$$\cdots C_{j+1} \xrightarrow{\partial_{j+1}} C_j \xrightarrow{\partial_j} C_{j-1} \cdots$$

One can finally define the *image* and *kernel* of the boundary operator given by

$$\text{Im } \partial_j = \{c \in C_{j-1} \mid \exists c' \in C_j \text{ where } \partial_j c' = c\}$$

and

$$\text{Ker } \partial_j = \{c \in C_j \mid \partial_j c = 0\}.$$

Now $\text{Im } \partial_{j+1}$ is the vector space of j -boundaries and $\text{Ker } \partial_j$ is the vector space of j cycles. Then the j th homology group is the quotient of $\text{Ker } \partial_j$ over $\text{Im } \partial_{j+1}$, i.e.,

$$H_j = \text{Ker } \partial_j / \text{Im } \partial_{j+1}.$$

The calculation of H_j is seen to be an exercise in linear algebra given the matrix representation of the boundary operators. In our exposition here, we have been working over a field resulting in vector spaces which are characterized purely by their ranks (the Betti numbers). One approach to this is also via the combinatorial Laplacian as outlined in [13]. More generally, one can work over a ring and H_j would then be an Abelian group.

4. The Deformation Retract Argument

In this section we prove Proposition 3.1. Recall that $\varepsilon < \sqrt{3/5}\tau$. Consider the canonical map $\pi: U \rightarrow \mathcal{M}$ given by (π is the restriction of π_0 to U)

$$\pi(x) = \arg \min_{p \in \mathcal{M}} \|x - p\|.$$

Then we see that the fibers $\pi^{-1}(p)$ are given by $T_p^\perp \cap U \cap B_\tau(p)$. The intersection with $B_\tau(p)$ is necessary to eliminate distant regions of U that may intersect with T_p (because the manifold may curve around over great distances) but do not belong to the fiber. For example, for the standard circle in \mathbb{R}^2 , at any point p on the circle, T_p^\perp intersects the circle at two points. One of these is in $B_\tau(p)$ and the other is not. Therefore,

$$\pi^{-1}(p) = \bigcup_{x \in \bar{x}} B_\varepsilon(x) \cap T_p^\perp \cap B_\tau(p),$$

where T_p^\perp is the normal subspace at $p \in \mathcal{M}$ orthogonal to the tangent space T_p . Let us also define $st(p)$ as

$$st(p) = \bigcup_{\{x \in \bar{x}; x \in B_\varepsilon(p)\}} B_\varepsilon(x) \cap T_p^\perp \cap B_\tau(p).$$

It is immediately clear that

$$st(p) \subseteq \pi^{-1}(p).$$

Then the following simple proposition is true.

Proposition 4.1. *$st(p)$ is star shaped relative to p and therefore contracts to p .*

Proof. Consider arbitrary $v \in st(p)$. Then $v \in B_\varepsilon(x) \cap T_p^\perp$ for some $x \in \bar{x}$ such that $x \in B_\varepsilon(p)$. Since $x \in B_\varepsilon(p)$, we immediately have $p \in B_\varepsilon(x)$. Since v, p are both in $B_\varepsilon(x)$, by convexity of Euclidean balls, we have that the line segment $\bar{v}p$ joining v to p is entirely contained in $B_\varepsilon(x)$. At the same time, $\bar{v}p$ is entirely contained in T_p^\perp and it follows therefore that $\bar{v}p$ is contained in $st(p)$. \square

We next show that the inclusion of $st(p)$ in $\pi^{-1}(p)$ is an equality proving that $\pi^{-1}(p)$ contracts to p .

Proposition 4.2.

$$st(p) = \pi^{-1}(p).$$

Proof. We need to show that $\pi^{-1}(p) \subseteq st(p)$. Consider an arbitrary $v \in B_\varepsilon(q) \cap T_p^\perp \cap B_\tau(p)$ where $q \in \bar{x}$ and $q \notin B_\varepsilon(p)$. For such v the picture of Fig. 1 can be drawn. Following Lemma 4.1, we see that the distance of v to p is at most ε^2/τ . Now by the fact that \bar{x} is $(\varepsilon/2)$ -dense, we have that there is some point $x \in \bar{x}$ which is within $\varepsilon/2$ of p . The worst-case picture of this is shown in Fig. 2. From Lemma 4.2, we see that $v \in B_\varepsilon(x)$ for this x . The proposition is proved. \square

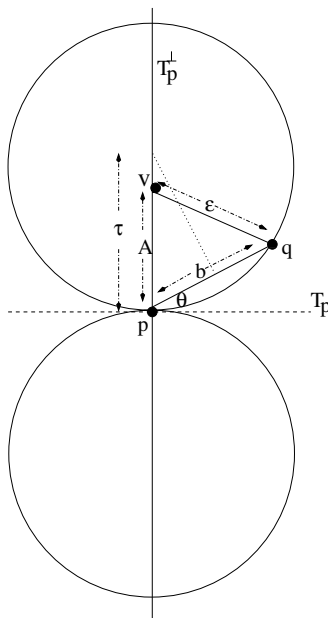


Fig. 1. A picture showing the worst case. The picture shows the plane passing through points v , p , and q . T_p and T_p^\perp are shown intersecting with this plane and are represented by the dotted horizontal line and the solid vertical line, respectively. On the plane of interest, one may then draw two circles (of radius τ each) that are tangent to T_p and are on either side of T_p as shown. Clearly, v lies on T_p^\perp and is marked in the figure. On the other hand, q could potentially lie anywhere outside the two circles. A moment's reflection shows that $\|v - p\|$ is greatest when q lies on one of the two circles. Without loss of generality one may consider it to lie on the top circle as shown. Over all choices of such q , the worst case is derived in Lemma 4.1

These two propositions taken together show that \mathcal{M} is a deformation retract of U . We see that $\mathcal{M} \subset U$. Further let $F(x, t): U \times [0, 1] \rightarrow U$ be given by $F(x, t) = tx + (1 - t)\pi(x)$. Then F is continuous, $F(x, 0) = \pi$, and $F(x, 1)$ is the identity map.

Lemma 4.1. *Consider any $q \notin B_\varepsilon(p)$. Let $v \in B_\varepsilon(q) \cap T_p^\perp \cap B_\tau(p)$. Then the Euclidean distance from v to p is less than ε^2/τ .*

Proof. We need to consider which configuration of v , q , and p makes the distance $\|v - p\|$ as large as possible. It is easiest to reason about this in the plane passing through these points. It suffices to consider q on the curve as shown in Fig. 1. See the caption for further explanation. Following the symbols on the figure, we have

$$A = b \sin(\theta) + \sqrt{\varepsilon^2 - b^2 \cos^2(\theta)},$$

where $b = 2\tau \sin(\theta)$. Therefore, we have

$$A = 2\tau \sin^2(\theta) + \sqrt{\varepsilon^2 - 4\tau^2 \sin^2(\theta) \cos^2(\theta)}.$$

From this we see that

$$\begin{aligned} \frac{dA}{d\theta} &= 2\tau \sin(2\theta) - \frac{4\tau^2 \sin(2\theta) \cos(2\theta)}{2\sqrt{\varepsilon^2 - \tau^2 \sin^2(2\theta)}} \\ &= 2\tau \sin(2\theta) \left(1 - \frac{\tau \cos(2\theta)}{\sqrt{\varepsilon^2 - \tau^2 \sin^2(2\theta)}} \right). \end{aligned}$$

It is easy to check that if $\varepsilon < \tau$, $dA/d\theta < 0$, i.e., A is monotonically decreasing with θ . Therefore the worst-case situation is when $b = 2\tau \sin(\theta) = \varepsilon$. For this value of θ , we see that $A = \varepsilon^2/\tau$. \square

The following lemma ensures that there is an $x \in \bar{x} \cap B_\varepsilon(p)$ such that $v \in B_\varepsilon(x) \cap T_p^\perp$.

Lemma 4.2. *Let \bar{x} be $(\varepsilon/2)$ -dense in \mathcal{M} . For any $p \in \mathcal{M}$, let $v \in \pi^{-1}(p)$. Then for $0 < \varepsilon < \sqrt{3/5}\tau$, we have that $v \in B_\varepsilon(x) \cap T_p^\perp$ for some $x \in B_\varepsilon(p) \cap \bar{x}$.*

Proof. By the $(\varepsilon/2)$ -dense property, we know that there is an $x \in \bar{x}$ such that $x \in B_{\varepsilon/2}(p)$. Consider the picture in Fig. 2. This represents the most unfavorable position that such an x might have for the current context. The picture shows the plane passing through the points x , v , and p . By the same argument of Lemma 4.1 we see that

$$A = \sqrt{\varepsilon^2 - b^2 \cos^2(\theta)} - b \sin(\theta),$$

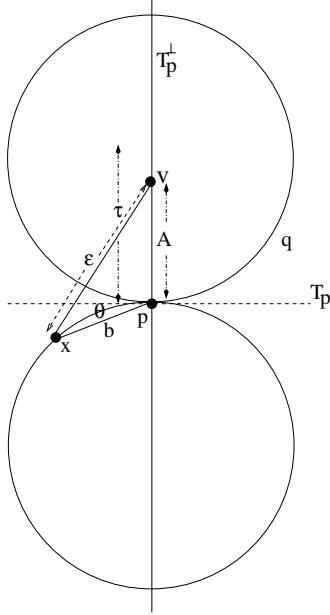


Fig. 2. A picture showing the worst case. The picture is of the plane containing the points p , v , and x . The two circles are each of radius τ and tangent to T_p . T_p and T_p^\perp are represented by their intersection with the plane of interest as dotted horizontal and solid vertical lines, respectively.

where $b = 2\tau \sin(\theta) = \varepsilon/2$. Putting this value in, we have

$$A = \sqrt{\varepsilon^2 - \frac{\varepsilon^2}{4} \left(1 - \frac{\varepsilon^2}{16\tau^2}\right)} - 2\tau \frac{\varepsilon^2}{16\tau^2}.$$

Simplifying, we see that $A > \varepsilon^2/\tau$ (needed by Lemma 4.1) if

$$\sqrt{\varepsilon^2 - \frac{\varepsilon^2}{4} \left(1 - \frac{\varepsilon^2}{16\tau^2}\right)} > \frac{9}{8} \frac{\varepsilon^2}{\tau}.$$

Squaring both sides, we have

$$\frac{3}{4}\varepsilon^2 + \frac{\varepsilon^4}{64\tau^2} > \frac{81\varepsilon^4}{64\tau^2}.$$

This simplifies to

$$\frac{\varepsilon^2}{\tau^2} < \frac{3}{5}.$$

Therefore, as long as $\varepsilon < \sqrt{\frac{3}{5}}\tau$, we will have that $v \in B_\varepsilon(x)$ for a suitable x . \square

5. Probability Bounds

Following our assumption, that the points x_i are drawn at random, we now provide a bound on how many examples need to be drawn so that the empirically constructed complex has the same homology as the manifold. We begin with a basic probability lemma.

Lemma 5.1. *Let $\{A_i\}$ for $i = 1, \dots, l$ be a finite collection of measurable sets and let μ be a probability measure on $\bigcup_{i=1}^l A_i$ such that for all $1 \leq i \leq l$, we have $\mu(A_i) > \alpha$. Let $\bar{x} = \{x_1, \dots, x_n\}$ be a set of n i.i.d. draws according to μ . Then if*

$$n \geq \frac{1}{\alpha} \left(\log l + \log \left(\frac{1}{\delta} \right) \right)$$

we are guaranteed that with probability $> 1 - \delta$, the following is true:

$$\forall i, \quad \bar{x} \cap A_i \neq \emptyset.$$

Proof. This follows from a simple application of the union bound. Let E_i be the event that $\bar{x} \cap A_i$ is empty. The probability with which this happens is given by

$$\mathbb{P}E_i = (1 - \mu(A_i))^n \leq (1 - \alpha)^n.$$

Therefore, by the union bound, we have

$$\mathbb{P}\bigcup_{i=1}^l E_i \leq \sum_{i=1}^l \mathbb{P}E_i \leq l(1 - \alpha)^n.$$

It remains to show that for $n \geq (1/\alpha)(\log l + \log(1/\delta))$, we have

$$l(1 - \alpha)^n \leq \delta.$$

To see this, simply note that $f(x) = xe^x - e^x + 1 \geq 0$ for all $x \geq 0$. This is seen by noting that $f(0) = 0$ and $f'(x) = xe^x \geq 0$ for all $x \geq 0$. Putting $x = \alpha$ in the above function, we have

$$(1 - \alpha) \leq e^{-\alpha}$$

and therefore it is easily seen that

$$l(1 - \alpha)^n \leq le^{-n\alpha} \leq \delta$$

for the appropriate choice of n . \square

Applying this to our setting, we consider a cover of the manifold \mathcal{M} by balls of radius $\varepsilon/4$. Let $\{y_i; 1 \leq i \leq l\}$ be the centers of such balls that constitute a minimal cover. Therefore, we can choose $A_i = B_{\varepsilon/4}(y_i) \cap \mathcal{M}$. Applying the above lemma, we immediately have an estimate on the number of examples we need to collect. This is given by

$$\frac{1}{\alpha} \left(\log l + \log \left(\frac{1}{\delta} \right) \right),$$

where

$$\alpha = \min_i \frac{\text{vol}(A_i)}{\text{vol}(\mathcal{M})}$$

and l is the $\varepsilon/4$ covering number. These may be expressed entirely in terms of natural invariants of the manifold and we derive these quantities below.

First, we note that the covering number may be bounded in terms of the packing number, i.e., the maximum number of sets of the form $N_i = B_r \cap \mathcal{M}$ (at scale r) that may be packed into \mathcal{M} without overlap. In particular, if $C(\varepsilon)$ is the ε -covering number of \mathcal{M} and $P(\varepsilon)$ is the ε -packing number, then the following simple lemma is true.

Lemma 5.2.

$$P(2\varepsilon) \leq C(2\varepsilon) \leq P(\varepsilon).$$

Proof. The fact that $P(2\varepsilon) \leq C(2\varepsilon)$ follows from the definition. To see that $C(2\varepsilon) \leq P(\varepsilon)$, begin by letting $B_\varepsilon(x_1), \dots, B_\varepsilon(x_N)$ be a realization of an optimal ε -packing so that $N = P(\varepsilon)$. We claim that $B_{2\varepsilon}(x_1), \dots, B_{2\varepsilon}(x_N)$ form a 2ε -cover. If not, there exists an $x \in \mathcal{M}$ such that $B_\varepsilon(x) \cap B_\varepsilon(x_i)$ is empty for all i . In that case, one can add $B_\varepsilon(x)$ to the collection to increase the packing number by 1 leading to a contradiction. Since $B_{2\varepsilon}(x_1), \dots, B_{2\varepsilon}(x_N)$ is a valid 2ε -cover, we have $C(2\varepsilon) \leq N = P(\varepsilon)$. \square

Since l is the $\varepsilon/4$ covering number, we see that $l \leq P(\varepsilon/8)$ from Lemma 5.2. Now we need to bound the packing number. To do so, we need the following result.

Lemma 5.3. *Let $p \in \mathcal{M}$. Now consider $A = \mathcal{M} \cap B_\varepsilon(p)$. Then $\text{vol}(A) \geq (\cos(\theta))^k \text{vol}(B_\varepsilon^k(p))$ where $B_\varepsilon^k(p)$ is the k -dimensional ball in T_p centered at p , $\theta = \arcsin(\varepsilon/2\tau)$. All volumes are k -dimensional volumes where k is the dimension of \mathcal{M} .*

Proof. Consider the tangent space at p given by T_p and let f be the projection of \mathbb{R}^N to T_p . Let $B_r^k(p)$ be the k -dimensional ball of radius $r = \varepsilon \cos(\theta)$ (where $\theta = \arcsin(\varepsilon/2\tau)$) centered at p lying in T_p . Let $f_A = \{f(q) \mid q \in A\}$ be the image of A under f . We will show that $B_r^k(p) \subset f_A$. Since f is a projection we have

$$\text{vol}(A) \geq \text{vol}(f_A) \geq \text{vol}(B_r^k(p)) = (\cos^k(\theta)) \text{vol}(B_\varepsilon^k(p)).$$

To see that $B_r^k(p) \subset f_A$, notice that f is an open map whose derivative is non-singular for all $q \in A$ (by Lemma 5.4). Therefore f is locally invertible and there exists a ball $B_s^k(p)$ of radius s such that $f^{-1}(B_s^k(p)) \subset A$. One can keep increasing s until it happens for the first time (say at $s = s'$) that $f^{-1}(B_{s'}^k(p)) \not\subset A$. At this stage, there exists a point q in the closure of A such that either (i) f is singular at q or (ii) $q \notin A$. By Lemma 5.4, we see that (i) is impossible. Therefore, $q \notin A$ but q is in the closure of A implying that $\|q - p\| = \varepsilon$. We see that $s' = \varepsilon \cos(\varphi)$ where φ is the angle between the line $\bar{q}p$ (the line joining q to p) and the line $f(q)p$ (the line joining $f(q)$ to p). By the curvature bound implied by τ , we see that $|\varphi| \leq |\theta|$ and therefore $s' = \varepsilon \cos(\varphi) \geq \varepsilon \cos(\theta) = r$. \square

Lemma 5.4. *Let $p \in \mathcal{M}$, let $A = \mathcal{M} \cap B_\varepsilon(p)$, and let f be the projection to the tangent space at p (T_p). Then for all $\varepsilon < \tau/2$, the derivative df is non-singular at all points $q \in A$.*

Proof. Suppose df was singular for some $q \in A$. That means that the tangent space at q (T_q) is oriented so that the vector with origin q and endpoint $f(q)$ lies in T_q . Since $q \in B_\varepsilon(p)$, we have that $d = \|q - p\| < \tau/2$. Putting Propositions 6.2 and 6.3 together, we get that

$$\cos(\varphi) \geq \sqrt{1 - \frac{2d}{\tau}} > 0,$$

where φ is the angle between T_p and T_q . From this we see that $\varphi < \pi/2$ leading to a contradiction. \square

Using Lemma 5.3, we see that a simple bound on the packing number is obtained. We obtain immediately that

$$P(\varepsilon) \leq \frac{\text{vol}(\mathcal{M})}{(\cos^k(\theta)) \text{vol}(B_\varepsilon^k(p))}.$$

Therefore, we have

$$l \leq P\left(\frac{\varepsilon}{8}\right) \leq \frac{\text{vol}(\mathcal{M})}{(\cos^k(\theta_2)) \text{vol}(B_{\frac{\varepsilon}{8}}^k(p))},$$

where $\theta_2 = \arcsin(\varepsilon/16\tau)$. Similarly, we have that

$$\frac{1}{\alpha} \leq \frac{\text{vol}(\mathcal{M})}{(\cos^k(\theta_1)) \text{vol}(B_{\frac{\varepsilon}{4}}^k(p))},$$

where $\theta_1 = \arcsin(\varepsilon/8\tau)$.

6. Curvature and the Condition Number $1/\tau$

In this section¹ we examine the consequences of the condition number $1/\tau$ for the submanifold \mathcal{M} . As we have mentioned before, τ controls the curvature of the manifold at every point. This fact has been exploited in our earlier proofs. For submanifolds, one may formally study curvature through the second fundamental form (see, e.g., [9]). Here we show formally that the norm of the second fundamental form is bounded by $1/\tau$. Thus a large τ corresponds to a well-conditioned submanifold that has low curvature.

Proposition 6.1 states the bound on the norm of the second fundamental form. Proposition 6.2 states a bound on the maximum angle between tangent spaces at different points in \mathcal{M} . Proposition 6.3 states a bound on the maximum difference between the geodesic distance and the ambient distance for neighboring points in \mathcal{M} .

We begin by recalling the second fundamental form. Fix a point $p \in \mathcal{M}$. Following standard accounts (see, e.g., [9]), there exists a symmetric bilinear form $B: T_p \times T_p \rightarrow T_p^\perp$ that maps any two vectors in the tangent space ($u, v \in T_p$) into a vector $B(u, v)$ in the normal space. Thus for any normal vector (unit norm) $\eta \in T_p^\perp$, one can define the following:

$$B_\eta(u, v) = \langle \eta, B(u, v) \rangle = \langle u, L_\eta v \rangle,$$

where the inner product $\langle \cdot, \cdot \rangle$ is the usual inner product in the tangent space of the ambient manifold (in our case \mathbb{R}^N). Since $B_\eta: T_p \times T_p \rightarrow \mathbb{R}$ is symmetric and bilinear, we see that $L_\eta: T_p \rightarrow T_p$ is a linear self-adjoint operator. The norm of the second fundamental form in direction η is now given by

$$\lambda_\eta = \sup_{u \in T_p} \frac{\langle u, L_\eta u \rangle}{\langle u, u \rangle}.$$

It is seen that λ_η is the largest eigenvalue of L_η . (In general, the eigenvalues are also known as the principal curvatures in the normal direction η .) Given this, we can prove the following proposition that characterizes the relation between the curvature through the second fundamental form and the condition number of the submanifold.

Proposition 6.1. *If \mathcal{M} is a submanifold of \mathbb{R}^N with condition number $1/\tau$, then the norm of the second fundamental form is bounded by $1/\tau$ in all directions. In other words, for all points $p \in \mathcal{M}$ and for all (unit norm) $\eta \in T_p^\perp$, we have*

$$\lambda_\eta = \sup_{u \in T_p} \frac{\langle u, L_\eta u \rangle}{\langle u, u \rangle} \leq \frac{1}{\tau}.$$

¹ Thanks to Nat Smale for discussions leading to the writing of this section.

Proof. We prove by contradiction. Suppose the proposition is false. Then there exists a point $p \in \mathcal{M}$, a tangent vector (unit norm) $u \in T_p$, and a normal vector (unit norm) η such that

$$\langle \eta, B(u, u) \rangle > \frac{1}{\tau}.$$

Consider a geodesic curve $c(t) \in \mathcal{M}$ parametrized by arc length such that $c(0) = p$ and $\dot{c}(0) = (dc/dt)(0) = u$. For convenience, we place the origin at p so that $c(0) = 0 = p$. With this (ambient) coordinate system, consider the point given by $\tau\eta$, i.e., the point a distance τ from p in the direction η . By our hypothesis on the condition number of the submanifold, we see that $p \in \mathcal{M}$ is the closest point on the manifold to the center of the τ -ball given by $\tau\eta$:

$$\text{for all } t, \quad \|c(t) - \tau\eta\|^2 \geq \tau^2$$

from which we get

$$\text{for all } t, \quad \langle c(t), c(t) \rangle - 2\tau \langle c(t), \eta \rangle \geq 0.$$

Consider the function $g(t) = \langle c(t), c(t) \rangle - 2\tau \langle c(t), \eta \rangle$. Since $c(0) = 0$, we see that $g(0) = 0$. Further, we have $g'(t) = 2\langle c(t), \dot{c}(t) \rangle - 2\tau \langle \dot{c}(t), \eta \rangle$. Since $c(0) = 0$ and $\langle \dot{c}(0), \eta \rangle = 0$, we see that $g'(0) = 0$. Finally, $g''(t) = 2\langle \dot{c}(t), \dot{c}(t) \rangle + 2\langle c(t), \ddot{c}(t) \rangle - 2\tau \langle \ddot{c}(t), \eta \rangle$. Since c is parametrized by arc length, we have $\langle \dot{c}(t), \dot{c}(t) \rangle = 1$ and $g''(0) = 2 - 2\tau \langle \ddot{c}(0), \eta \rangle$.

Noting that the tangent vector field dc/dt is parallel (see the proof of Proposition 6.2), we see that $B(dc/dt, dc/dt) = \ddot{c}(t)$. Therefore, by assumption, we have that

$$\langle \eta, B(u, u) \rangle = \left\langle \eta, B\left(\frac{dc}{dt}, \frac{dc}{dt}\right) \right\rangle = \langle \eta, \ddot{c}(0) \rangle > \frac{1}{\tau}.$$

Therefore, $g''(0) < 2 - 2\tau(1/\tau) = 0$. By continuity, there exists a t^* such that $g(t^*) < 0$. However, this leads to a contradiction since $g(t) \geq 0$ for all t . \square

Since the norm of the second fundamental form is bounded, we see that the manifold cannot curve too much locally. As a result, the angle between tangent spaces at nearby points cannot be too large. Let p and q be two points in the submanifold \mathcal{M} with associated tangent spaces T_p and T_q . Since T_p and T_q are affine subspaces of \mathbb{R}^N , one can compare them in the ambient space in a standard way.

Formally, one may transport the tangent spaces to the origin (according to the standard connection defined in the ambient space \mathbb{R}^N) and then compare vectors in each of these tangent spaces with each other. Thus for any (unit norm) vectors $u \in T_p$ and $v \in T_q$, we may define the angle θ between them by

$$\cos(\theta) = |\langle u', v' \rangle|,$$

where $\langle \cdot, \cdot \rangle$ is the usual inner product in \mathbb{R}^N , and u', v' are the vectors obtained by parallel transport (in \mathbb{R}^N) of u and v respectively, to the origin. Hereafter, we always take this construction as standard. We drop the prime notation and use $\langle u, v \rangle$ to denote $\langle u', v' \rangle$ in what follows.

We can now state the following proposition.

Proposition 6.2. *Let \mathcal{M} be a submanifold of \mathbb{R}^N with condition number $1/\tau$. Let $p, q \in \mathcal{M}$ be two points with geodesic distance given by $d_{\mathcal{M}}(p, q)$. Let φ be the angle between the tangent spaces T_p and T_q defined by $\cos(\varphi) = \min_{u \in T_p} \max_{v \in T_q} |\langle u, v \rangle|$. Then $\cos(\varphi)$ is greater than $1 - (1/\tau)d_{\mathcal{M}}(p, q)$.*

Proof. Consider two points $p, q \in \mathcal{M}$ connected by a geodesic curve $c(t) \in \mathcal{M}$. Let $c(t)$ be parametrized (proportional to arc length) so that $c(0) = p$, and $c(1) = q$.

Now let $v_p \in T_p$ be a tangent vector (unit norm) and let $v(t)$ be the parallel transport of this vector along the curve $c(t)$. Thus we have $v(0) = v_p$, $v(1) = v_q \in T_q$. Clearly, $\langle v(t), v(t) \rangle = 1$ for all t since v is parallel.

Notice that

$$\langle v(0), v(1) \rangle = \langle v(0), v(0) + w \rangle = 1 + \langle v(0), w \rangle, \quad (1)$$

where

$$w = \int_0^1 \left(\frac{dv}{dt} \right) dt. \quad (2)$$

Combining (1) and (2), we see

$$\cos(\theta) = |\langle v(0), v(1) \rangle| \geq 1 - |\langle v(0), w \rangle| \geq 1 - \|w\|, \quad (3)$$

where θ is the angle between the vectors $v(0)$ and $v(1)$. Since $v_p = v(0)$ was arbitrary, it is easy to check that $\cos(\varphi) \geq \cos(\theta)$.

Now

$$\frac{dv}{dt} = \bar{\nabla}_{dc/dt} v(t),$$

where $\bar{\nabla}$ denotes the connection in Euclidean space. At the same time

$$\nabla_{dc/dt} v(t) = (\bar{\nabla}_{dc/dt} v(t))^T,$$

where for any $r \in \mathcal{M}$ and $v \in \bar{T}_r$ (here \bar{T}_r is the tangent space of \mathbb{R}^N at r) we denote by $(v)^T$ the projection of v onto T_r (here T_r is the tangent space to \mathcal{M} at r viewed as an affine space with origin r). However, since $v(t)$ is parallel, we have that $\nabla_{dc/dt} v(t) = 0$. Therefore, $\bar{\nabla}_{dc/dt} v(t)$ is entirely in the space normal to $T_{c(t)}$, but the component of $\bar{\nabla}_{dc/dt} v(t)$ in the normal direction is precisely given by the second fundamental form. Hence, we have that

$$\frac{dv}{dt} = B \left(\frac{dc}{dt}, v(t) \right),$$

where B is a symmetric, bilinear form (the second fundamental form). Letting η be a unit norm vector in the direction dv/dt , i.e., $\eta = (1/\|dv/dt\|)(dv/dt)$, we see that

$$\left\| \frac{dv}{dt} \right\| = \left\langle \eta, \frac{dv}{dt} \right\rangle = \left\langle \eta, B \left(\frac{dc}{dt}, v(t) \right) \right\rangle = \left\langle \frac{dc}{dt}, L_n v(t) \right\rangle,$$

where L_n is a self-adjoint linear operator. By Proposition 6.1, the norm of L_n is bounded by $1/\tau$. Therefore, we have

$$\left\| \frac{dv}{dt} \right\| \leq \left\| \frac{dc}{dt} \right\| \|L_n v\| \leq \left\| \frac{dc}{dt} \right\| \|L_n\|$$

and

$$\|w\| = \left\| \int_0^1 \frac{dv}{dt} \right\| \leq \int_0^1 \left\| \frac{dv}{dt} \right\| \leq \|L_n\| \int_0^1 \left\| \frac{dc}{dt} \right\| dt \leq \frac{1}{\tau} d_{\mathcal{M}}(p, q). \quad (4)$$

Combining (3) and (4), we get

$$\cos(\varphi) \geq 1 - \frac{1}{\tau} d_{\mathcal{M}}(p, q). \quad \square$$

We next show a relationship between the geodesic distance $d_{\mathcal{M}}(p, q)$ and the ambient distance $\|p - q\|_{\mathbb{R}^N}$ for any two points p and q on the submanifold \mathcal{M} .

Proposition 6.3. *Let \mathcal{M} be a submanifold of \mathbb{R}^N with condition number $1/\tau$. Let p and q be two points in \mathcal{M} such that $\|p - q\|_{\mathbb{R}^N} = d$. Then for all $d \leq \tau/2$, the geodesic distance $d_{\mathcal{M}}(p, q)$ is bounded by*

$$d_{\mathcal{M}}(p, q) \leq \tau - \tau \sqrt{1 - \frac{2d}{\tau}}.$$

Proof. Consider two points $p, q \in \mathcal{M}$ and let $c(t)$ be a geodesic curve joining them such that $c(0) = p$ and $c(s) = q$. Let c be parametrized by arc length so that $\|\dot{c}(t)\| = 1$ for all t and $d_{\mathcal{M}}(p, q) = s$.

Noting that the tangent vector field \dot{c} along the curve is parallel, we have $\ddot{c} = B(\dot{c}, \dot{c})$ and from Proposition 6.1 we see that for all t ,

$$\|\ddot{c}\| = \|B(\dot{c}, \dot{c})\| \leq \frac{1}{\tau}.$$

The chord length between p and q is given by $\|c(s) - c(0)\|$ and we now relate this to the geodesic distance $d_{\mathcal{M}}(p, q)$. Observe that

$$c(s) - c(0) = \int_0^s \dot{c}(t) dt.$$

Now

$$\dot{c}(t) = \dot{c}(0) + \int_0^t \ddot{c}(r) dr.$$

Thus $\dot{c}(t) = \dot{c}(0) + u(t)$ where $u(t) = \int_0^t \ddot{c}(r) dr$. We see that

$$\|u(t)\| \leq \int_0^t \|\ddot{c}(r)\| dr \leq \frac{t}{\tau}.$$

Therefore,

$$\|c(s) - c(0)\| = \left\| \int_0^s \dot{c}(0) dt + \int_0^s u(t) dt \right\| \geq s \|\dot{c}(0)\| - \int_0^s \|u(t)\| dt \geq s - \int_0^s \frac{t}{\tau} dt.$$

Therefore we get

$$\|c(s) - c(0)\| = d \geq s - \frac{s^2}{2\tau}, \quad (5)$$

where d is the ambient distance between the points p and q while s is the geodesic distance between these same points. The inequality in (5) is satisfied only if $s \leq \tau - \tau\sqrt{1 - 2d/\tau}$ or $s \geq \tau + \tau\sqrt{1 - 2d/\tau}$. Since $s = 0$ when $d = 0$, we know that the second inequality does not apply. Therefore, from the first inequality, we have

$$s \leq \tau - \tau\sqrt{1 - \frac{2d}{\tau}}. \quad \square$$

7. Handling Noisy Data

In this section we show that if our data is noisy in the sense that it is drawn from a probability distribution that is concentrated around (rather than on) the manifold, the homology of the manifold can still be computed from noisy data.

7.1. The Model of Noise

Consider a probability measure μ concentrated around the manifold. We assume that μ satisfies the following two regularity conditions:

1. The support of μ ($\text{supp } \mu$) is contained in the tubular neighborhood of radius r around \mathcal{M} . Thus $\text{supp } \mu \subset \text{Tub}_r(\mathcal{M})$.
2. For every $0 < s < r$, we have that

$$\inf_{p \in \mathcal{M}} \mu(B_s(p)) > k_s,$$

where k_s is a constant depending on s and independent of p .

In what follows we assume the data is drawn in i.i.d. fashion according to a P that satisfies the above properties.

7.2. Main Topological Lemma: Sufficient Conditions

We proceed by constructing ε -balls centered on our data points. If these data are s -dense on the manifold, then the homology of the union of these balls will equal that of the manifold \mathcal{M} even if the data is drawn from a noisy distribution. In order to see that this might be the case, we provide a simple argument. This argument works with non-optimal

choices of ε and s and later sections enter into the considerations of choosing better values for these parameters and therefore providing more natural complexity estimates.

Let $\bar{x} = \{x_1, \dots, x_n\}$ be a set of n points in the tubular neighborhood of radius r around \mathcal{M} . Let U be given by

$$U = \bigcup_{x \in \bar{x}} B_\varepsilon(x).$$

Proposition 7.1. *If \bar{x} is r -dense in \mathcal{M} then \mathcal{M} is a deformation retract of U for all $r < (\sqrt{9} - \sqrt{8})\tau$ and*

$$\varepsilon \in \left(\frac{(r + \tau) - \sqrt{r^2 + \tau^2 - 6\tau r}}{2}, \frac{(r + \tau) + \sqrt{r^2 + \tau^2 - 6\tau r}}{2} \right).$$

Proof. We show that for each $p \in \mathcal{M}$, it is the case that $\pi^{-1}(p)$ contracts to p . Consider a $v \in \pi^{-1}(p)$. Consider the line segment, $\bar{v}p$, joining v to p . We claim that this line segment is entirely contained in $\pi^{-1}(p)$. Clearly, if $v \in B_\varepsilon(x)$ for some $x \in \bar{x} \cap B_\varepsilon(p)$, this is immediate by the convexity of balls in Euclidean space. So we only need to consider the situation where $v \in B_\varepsilon(x)$ for some $x \notin \bar{x} \cap B_\varepsilon(p)$. So let $v \in B_\varepsilon(q) \cap T_p^\perp$. Let

$$u = \arg \min_{x \in \bar{v}p \cap B_\varepsilon(q)} \|x - p\|.$$

As long as $u \in B_\varepsilon(x)$ for some $x \in \bar{x} \cap B_\varepsilon(p)$, we see that the line segment $\bar{u}p$ is contained in $\pi^{-1}(p)$ and therefore v contracts to p .

Since we choose $r < \varepsilon$, we are guaranteed that there is an $x \in \bar{x} \cap B_r(p) \subset B_\varepsilon(p)$. The worst-case picture is shown in Fig. 3. Following the symbols of the figure, as long as

$$\tau - A < \varepsilon - r,$$

we have that v contracts to p . Thus we need

$$(\tau - (\varepsilon - r))^2 < A^2 = (\tau - r)^2 - \varepsilon^2. \quad (6)$$

Expanding the squares, this reduces to

$$\varepsilon^2 - \varepsilon(\tau + r) + 2\tau r < 0.$$

This is a quadratic in ε and is satisfied for

$$\varepsilon \in \left(\frac{(r + \tau) - \sqrt{r^2 + \tau^2 - 6\tau r}}{2}, \frac{(r + \tau) + \sqrt{r^2 + \tau^2 - 6\tau r}}{2} \right) \quad (7)$$

provided

$$r^2 - 6\tau r + \tau^2 > 0.$$

This, in turn, is a quadratic in r and it is easy to check that it is satisfied as long as

$$r < (3 - 2\sqrt{2})\tau = (\sqrt{9} - \sqrt{8})\tau. \quad (8)$$

Thus we see that for r, ε satisfying (7) and (8), we have that v contracts to p . \square

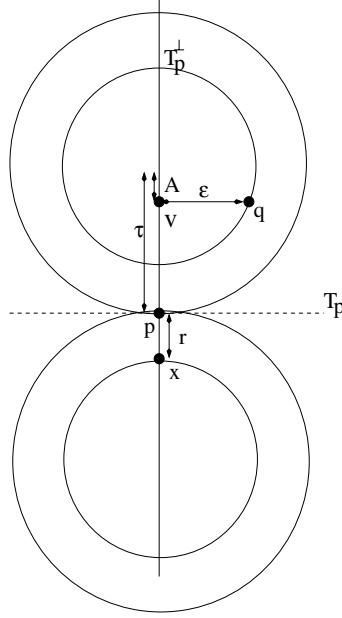


Fig. 3. A picture showing the worst case. As before, we draw the picture in the plane connecting points v , p , and q . T_p and T_p^\perp are intersected with this plane in the picture and shown by the dotted horizontal line and solid vertical line, respectively. The concentric circles have the same center and are of radius τ and $\tau - r$, respectively, and follow our usual construction in earlier figures and arguments. All lengths are marked by arrows.

We now need to compute the probability of drawing a random \bar{x} that is guaranteed to be r -dense. The following proposition is true.

Proposition 7.2. *Let $N_{r/2}$ be the $(r/2)$ -covering number of the manifold. Let $p_1, \dots, p_{N_{r/2}} \in \mathcal{M}$ be points on the manifold such that $B_{r/2}(p_i)$ realize an $(r/2)$ -cover of the manifold. Let \bar{x} be generated by i.i.d. draws according to a probability measure μ that satisfies the regularity properties described earlier. Then if $|\bar{x}| > (1/k_{r/2})(\log(N_{r/2}) + \log(1/\delta))$, with probability greater than $1 - \delta$, \bar{x} will be r -dense in \mathcal{M} .*

Proof. Take $A_i = B_{r/2}(p_i)$ and apply Lemma 5.1. By the conclusion of that lemma, we have that with high probability each of the A_i 's is occupied by at least one $x \in \bar{x}$. Therefore it follows that for any $p \in \mathcal{M}$, there is at least one $x \in \bar{x}$ such that $\|p - x\| < r$. Thus with high probability \bar{x} is r -dense on the manifold. \square

Putting these together, our main conclusion is

Theorem 7.1. *Let $N_{r/2}$ be the $(r/2)$ -covering number of the submanifold \mathcal{M} of \mathbb{R}^N . Let \bar{x} be generated by i.i.d. draws according to a probability measure μ that satisfies the regularity properties described earlier. Let $U = \bigcup_{x \in \bar{x}} B_\varepsilon(x)$. Then if $|\bar{x}| > (1/k_{r/2})(\log(N_{r/2}) + \log(1/\delta))$, with probability greater than $1 - \delta$, \mathcal{M} is a deformation*

retract of U as long as (i) $r < (\sqrt{9} - \sqrt{8})\tau$ and (ii)

$$\varepsilon \in \left(\frac{(r + \tau) - \sqrt{r^2 + \tau^2 - 6\tau r}}{2}, \frac{(r + \tau) + \sqrt{r^2 + \tau^2 - 6\tau r}}{2} \right).$$

7.3. *Main Topological Lemma—General Considerations*

In general, we may demand points that are s -dense. Putting ε -balls around these points we construct U in the usual way. The condition number τ and the noise bound r are additional parameters that are outside our control and determined externally. We now ask what is the feasible space $(s, \varepsilon, r, \tau)$ that will guarantee that U is homotopy equivalent to \mathcal{M} ?

Following our usual logic, we see that the worst-case situation is given by Fig. 4. An arbitrary $v \in B_\varepsilon(q) \cap T_p^\perp \cap B_\tau(p)$ will contract to p if

$$B_\varepsilon(q) \cap B_\varepsilon(x) \cap \bar{v}p \neq \varnothing.$$

This is the same as requiring

$$(\tau - w)^2 < (\tau - r)^2 - \varepsilon^2. \tag{9}$$

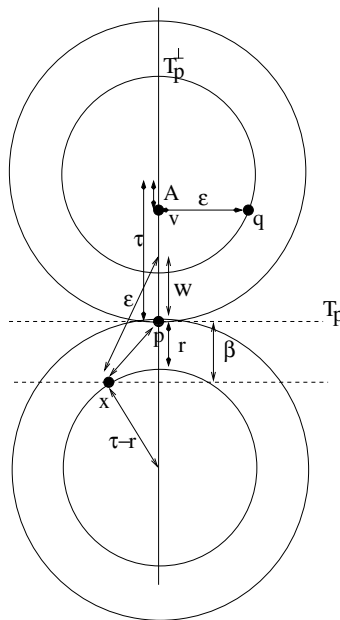


Fig. 4. A picture showing the worst case. As before, we draw the picture in the plane connecting points v , p , and q . T_p and T_p^\perp are intersected with this plane in the picture and shown by the dotted horizontal line and solid vertical line, respectively. The concentric circles have the same center and are of radius τ and $\tau - r$, respectively, and follow our usual construction in earlier figures and arguments. All lengths are marked by arrows.

Additionally, we have the following equations that need to be satisfied (following Fig. 4):

$$(\tau - r)^2 - (\tau - \beta)^2 = s^2 - \beta^2, \quad (10)$$

$$s^2 - \beta^2 + (\beta + w)^2 = \varepsilon^2. \quad (11)$$

If one eliminates w and β from the above equations, one will get a single inequality relating s, ε, τ, r that describes for each τ, r the feasible set of possible choices of s, ε that are sufficient to guarantee homotopy equivalence. Let us see how our earlier theorems follow from particular choices of this general set of equations.

7.3.1. The Case when $s = r$. We have already examined the case when the points \bar{x} are chosen to be r -dense in \mathcal{M} . Putting $s = r$ in (9)–(11), we see the following:

From (10), we have (for $s = r$)

$$(\tau - r)^2 - (\tau - \beta)^2 = r^2 - \beta^2.$$

This simplifies to give $\beta = r$.

Putting $\beta = r$ and $s = r$ in (11), we get

$$r^2 - r^2 + (r + w)^2 = \varepsilon^2,$$

giving us $w = \varepsilon - r$.

Finally, putting $w = \varepsilon - r$ in inequality (9), we get

$$(\tau - (\varepsilon - r))^2 < (\tau - r)^2 - \varepsilon^2,$$

which is the same as inequality (6) whose solution was examined in the previous section.

7.3.2. The Case when $r = 0$. We can recover our main theorem for the noise-free case by considering the case $r = 0$. We proceed to do this now.

The fundamental inequality of (9) gives us (for $r = 0$)

$$(\tau - w)^2 < \tau^2 - \varepsilon^2.$$

This is the same as requiring

$$w^2 - 2\tau + \varepsilon^2 < 0.$$

Using standard analysis for quadratic functions, we see that the following condition is required:

$$w > \tau - \sqrt{\tau^2 - \varepsilon^2}. \quad (12)$$

We can eliminate w using (10) and (11). Thus, from (10), we get $\beta = s^2/2\tau$ and substituting in (11), we get a quadratic equation in w whose positive solution is given by $w = -s^2/2\tau + \sqrt{s^4/4\tau^2 + (\varepsilon^2 - s^2)}$. This gives rise to the following condition:

$$-\frac{s^2}{2\tau} + \sqrt{\frac{s^4}{4\tau^2} + (\varepsilon^2 - s^2)} > \tau - \sqrt{\tau^2 - \varepsilon^2}. \quad (13)$$

Inequality (13) gives the feasible region for s and ε for the homotopy equivalence of U and \mathcal{M} . Let us consider the special case when $s = \varepsilon/2$ —a choice we made in Section 3 without any attention to optimality. Putting in this value, after several simplifying steps, one obtains that

$$\varepsilon^4 + 51\varepsilon^2\tau^2 - 48\tau^4 < 0. \quad (14)$$

This is satisfied for all $0 < \varepsilon^2 < 0.9244\tau^2$ or $0 < \varepsilon < 0.96\tau$.

Remark 1. Note that in our original proof of our main noise free theorem (Theorem 3.1), the deformation retract argument of Section 3 passes through the construction of $st(p)$ and shows contraction of $\pi^{-1}(p)$ by equating it with $st(p)$. This condition is stronger than we require. Here we see that the condition $B_\varepsilon(q) \cap B_\varepsilon(x) \cap \bar{v}p \neq \varnothing$ is sufficient. This latter condition is weaker and therefore gives us a slightly stronger version of Theorem 3.1 in the sense that it holds for a larger range of ε .

Remark 2. If we assume that τ, r are beyond our control, the sample complexity depends entirely upon s . Therefore if we wish to proceed by drawing the fewest number of examples, then it is necessary to maximize s subject to the condition of (13).

Remark 3. The total complexity of finding the homology depends both upon s and ε in a more complicated way. The size of \bar{x} depends entirely upon s and nothing else. However, the number of k -tuples to consider in the simplicial complex depends both upon the size of \bar{x} as well as ε because ε determines how many balls will have non-empty intersections. We leave this more nuanced complexity analysis for future consideration.

References

1. N. Amenta and M. Bern. Surface reconstruction by Voronoi filtering. *Discrete & Computational Geometry*, **22** (1999), 481–504.
2. N. Amenta, S. Choi, T. K. Dey, and N. Leekha. A simple algorithm for homeomorphic surface reconstruction. *International Journal of Computational Geometry Applications*, **12** (2002), 125–141.
3. M. Belkin and P. Niyogi. Semisupervised learning on Riemannian manifolds. *Machine Learning*, **56** (2004), 209–239.
4. A. Björner. Topological methods. In *Handbook of Combinatorics* (R. Graham, M. Grötschel, L. Lovász, eds.), pp. 1819–1872. North-Holland, Amsterdam, 1995.
5. F. Chazal and A. Lieutier. Weak feature size and persistent homology: computing homology of solids in \mathbb{R}^n from noisy data samples. Preprint.
6. S. W. Cheng, T. K. Dey, and E. A. Ramos. Manifold reconstruction from point samples. *Proceedings of ACM–SIAM Symposium on Discrete Algorithms*, pp. 1018–1027, 2005.
7. D. Cohen-Steiner, H. Edelsbrunner, and J. Harer. Stability of persistence diagrams, *Proceedings of the 21st Symposium on Computational Geometry*, pp. 263–271, 2005.
8. T. K. Dey, H. Edelsbrunner, and S. Guha. Computational topology. In *Advances in Discrete and Computational Geometry* (B. Chazelle, J. E. Goodman, and R. Pollack, eds.), pp. 109–143. Contemporary Mathematics 223. AMS, Providence, RI, 1999.
9. M. P. Do Carmo. *Riemannian Geometry*. Birkhäuser, Basel, 1992.
10. D. Donoho and C. Grimes. Hessian eigenmaps: new locally-linear embedding techniques for high-dimensional data. Preprint. Department of Statistics, Stanford University, 2003.
11. H. Edelsbrunner and E. P. Mücke. Three-dimensional alpha shapes. *ACM Transactions on Graphics*, **13** (1994), 43–72.

12. K. Fischer, B. Gaertner, and M. Kutz. Fast smallest-enclosing-ball computation in high dimensions. *Proceedings of the 11th Annual European Symposium on Algorithms (ESA)*, pp. 630–641, 2003.
13. J. Friedman. Computing Betti numbers via combinatorial laplacians. *Algorithmica*, **21** (1998), 331–346.
14. T. Kaczynski, K. Mischaikow, and M. Mrozek. *Computational Homology*. Springer-Verlag, New York, 2004.
15. J. Munkres. *Elements of Algebraic Topology*. Addison-Wesley, Menlo Park, CA, 1984.
16. S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, **290** (2000), 2323–2326.
17. J. B. Tenenbaum, V. De Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, **290** (2000), 2319–2323.
18. L. G. Valiant. A theory of the learnable. *Communications of the ACM*, **27**(11) (1984), 1134–1142.
19. Website for Smallest Enclosing Ball Algorithm. <http://www2.inf.ethz.ch/personal/gaertner/miniball.html>
20. A. Zomorodian and G. Carlsson. Computing persistent homology. *Discrete & Computational Geometry*, **33** (2005), 249–274.

Received June 21, 2005, and in revised form December 29, 2005, and March 16, 2006.

Online publication September 25, 2006.